# Automatic Synthesis of Parallel Unix Commands and Pipelines with KumQuat

Jiasi Shen
MIT EECS & CSAIL
USA
jiasi@csail.mit.edu

Martin Rinard
MIT EECS & CSAIL
USA
rinard@csail.mit.edu

Nikos Vasilakis
MIT EECS & CSAIL
USA
nikos@vasilak.is

## Abstract

We present KumQuat, a system for automatically generating data parallel implementations of Unix shell commands and pipelines. The generated parallel versions split input streams, execute multiple instantiations of the original pipeline commands to process the splits in parallel, then combine the resulting parallel outputs to produce the final output stream. KumQuat automatically synthesizes the combine operators, with a domain-specific combiner language acting as a strong regularizer that promotes efficient inference of correct combiners.

We evaluate KumQuat on 70 benchmark scripts that together have a total of 427 stages. KumQuat synthesizes a correct combiner for 113 of the 121 unique commands that appear in these benchmark scripts. The synthesis times vary between 39 seconds and 331 seconds with a median of 60 seconds. We present experimental results that show that these combiners enable the effective parallelization of our benchmark scripts.

## 1 Introduction

The Unix shell, working in tandem with the wide range of commands it supports, provides a convenient programming environment for many stream processing computations. Shell commands—which can be written in multiple languages—typically execute sequentially on a single processor. This sequential execution often leaves available data parallelism, in which a command operates on different parts of an input stream in parallel, unexploited. This observation has motivated the development of systems that exploit data parallelism in shell pipelines [18, 26]. A key prerequisite is obtaining the combiners required to merge the resulting multiple parallel output streams correctly into a single output stream. Previous systems rely on developers to manually implement such combiners and associate them with their corresponding shell commands [18, 26].

We present a new system, KumQuat, for automatically exploiting data parallelism available in Unix pipelines. Working with the commands in the pipeline as black boxes, KumQuat automatically generates inputs that explore the behavior of the command to infer and automatically generate a combiner for the command. This capability enables KumQuat to automatically generate data parallel versions of Unix pipelines, including pipelines that contain new commands or command options for which combiners were previously unavailable.

KumQuat targets commands that can be expressed as data parallel divide-and-conquer computations with two phases:[1] the first phase executes the original, unmodified command in parallel on disjoint parts of the input; the second phase combines the partial results from the first phase to obtain the final output. KumQuat generates candidate combiners, then repeatedly feeds selected inputs to parallelized versions of the command that use the candidate combiners. A comparison of the resulting outputs with corresponding outputs from the original serial version of the command enables KumQuat to identify a correct combiner for the command. A domain-specific combiner language acts as a strong regularizer that promotes efficient learning of correct combiners. The resulting (automatically generated) parallel computation executes directly in the same environment and with the same program and data locations as the original sequential command.

This paper makes the following contributions:

- **Algorithm:** We present a new algorithm that automatically synthesizes combiners for parallel and distributed versions of Unix commands. The resulting synthesized combiners enable the automatic generation of parallel and distributed versions of Unix commands and pipelines.

- **Domain-Specific Language:** We present a domain-specific language for combiner operators. This language supports both the class of combiner operators relevant to this domain and an efficient algorithm for automatically synthesizing these combiner operators.

- **Correct Combiners:** We present theorems (Theorem 2 and Theorem 4) that characterize when the combiner synthesis algorithm will identify a correct combiner for a given set of parallel output streams.

We also present an analysis of the interaction between our input generation algorithm, our benchmark commands, and our combiner synthesis algorithm. This analysis identifies why KumQuat generates correct combiners for the

---

[1]There is no requirement that the actual internal implementation must be structured as a divide-and-conquer computation — because KumQuat interacts with the command as a black box, the requirement is instead only that the computation that it implements can be expressed in this way.

113 of 121 benchmark commands for which correct combiners exist, identifies command patterns that ensure correct combiner synthesis and correct data parallel execution, and provides insight into the rationale behind the KumQuat design and the reasons why the KumQuat design can effectively exploit data parallelism available in its target class of commands.

- **Experimental Results:** We present experimental results that characterize the effectiveness of KumQuat on a set of 70 benchmark scripts that together have a total of 477 commands, among which 121 are unique and process an input stream. The results show that KumQuat can effectively synthesize combiners for the majority of our benchmark commands and that these synthesized combiners enable effective parallelizations of our benchmark Unix pipelines.

## 2 Example

Figure 1 presents an example pipeline that we use to illustrate KumQuat. The pipeline implements a computation that counts the frequency of words in an input document. The six commands in the pipeline (1) read the input document, (2) break the document into lines of words, (3) translate the words into lower case, (4) sort the words, (5) remove duplicates and prepend each unique word with a count, and (6) sort the words on their counts in reverse order.

The pipeline conforms to a standard Unix model that structures computations as pipelines of building-block commands that process character streams. The example commands process the streams as lines of words, with the lines and words separated by delimiters, in the example newline and space. As the commands process lines, words, or characters, they apply a function to each unit and either output the result of the function (commands "`tr -cs A-Za-z '\n'`" and "`tr A-Z a-z`"), sort the units according to a certain order (commands "`sort`" and "`sort -rn`"), or accumulate a result that is output when the command finishes reading the input stream and terminates (command "`uniq -c`").

To exploit the data parallelism available in this computation, KumQuat splits the input data stream into substreams, then instantiates the commands to process the input substreams in parallel. The result is a set of parallel output substreams that must then be combined to obtain the single final output stream of the command.

Different commands often require different combine operators. The combine operator for command "`tr A-Z a-z`" simply concatenates the output substreams. The combine operator for command "`tr -cs A-Za-z '\n'`" concatenates the output substreams, then reruns the command on this concatenated stream.[2] Note that `tr` commands with different flags may have different combine operators. The combine operators for sort commands apply an appropriate merge

---

[2]Simple concatenation is incorrect due to potential empty lines at the split boundary.

```
cat $IN | tr -cs A-Za-z '\n' | tr A-Z a-z |
    sort | uniq -c | sort -rn
```

**Figure 1.** Example pipeline that computes word frequencies [2, 26].

function, which may depend on the sort flag that specifies the comparison function. The "`uniq -c`" command produces a stream of (count, word) pairs. Given two streams $y_1$ and $y_2$, the combine operator compares the word in the last line of $y_1$ with the word in the first line of $y_2$. If they are the same, it concatenates $y_1$ and $y_2$ but combines the last and first lines to include the sum of the two word counts. Otherwise it simply concatenates $y_1$ and $y_2$. As these examples highlight, selecting an appropriate combine operator for each command is a critical step for obtaining a correct parallel execution.

The default KumQuat parallel computation applies the combine operator after the parallel execution of each command to obtain a single output stream for that command. In many cases, however, it is possible to enhance parallel performance by eliminating intermediate combine operators so that the output substreams for one command feed directly into the input substreams for the direct parallel execution of the next command. KumQuat therefore applies an optimization that automatically eliminates intermediate combine operators when possible (Section 3.5).

**Model of Computation**: KumQuat targets commands $f$ that have a combine operator $g$ that satisfies

$$f(x_1 ++ x_2) = g(f(x_1), f(x_2))$$

for all input streams $x_1, x_2$, where the streams are (potentially recursively) structured as units separated by delimiters. KumQuat currently targets character streams structured as lines with the newline delimiter, so that $x_1$ and $x_2$ terminate with newlines. ++ denotes string concatenation. A key step in the parallelization of $f$ is the synthesis of a correct combiner $g$ for $f$. To focus the synthesis algorithm on a productive space of candidate combiners, KumQuat works with combiners expressible in a domain-specific combiner language (Figure 3).

**Combiner Synthesis**: To infer a combiner $g$ for a command $f$, KumQuat works with a set of candidate combiner functions. In the current KumQuat implementation, this set consists of all combiner functions with seven or fewer nodes in the DSL abstract syntax tree. KumQuat repeatedly generates input streams $x_1$ and $x_2$, feeds the input streams to the serial and parallel versions of $f$ instantiated with candidate combiners $g$, and compares the resulting serial and parallel outputs to discard candidate combiners $g$ that do not satisfy the equation $f(x_1 ++ x_2) = g(f(x_1), f(x_2))$. The input stream generation algorithm is designed to produce inputs that quickly find and discard incorrect combiners (see below). With these generated inputs, we have found that the

combiner synthesis algorithm typically converges quickly to a few semantically equivalent correct combiners (Section 4).

**Input Generation**: KumQuat uses a set of *input shapes* to specify the format of each generated input stream. An input shape specifies the number of lines in the input stream, the number of words per line, and the number of characters per word. The input shape also specifies how diverse these input units are, in terms of the percentage of distinct lines, words, and characters. These input shapes are designed to generate meaningful inputs for commands that conform to our model of computation. A goal is to efficiently generate counterexample inputs that cause the command to produce counterexample outputs that enable KumQuat to identify and discard incorrect candidate combiners.

The design of input shapes is inspired by the observation that certain input shapes cause commands to produce outputs that enable KumQuat to identify and discard incorrect combiners. For example, when $f = ($tr -cs A-Za-z '\n'$)$ and $g = $ **concat**, a counterexample input has $x_1$ ending with a newline and $x_2$ starting with a newline. In this case, $f(x_1)$ also ends with a newline and $f(x_2)$ also starts with a newline, so $g(f(x_1), f(x_2))$ has two consecutive newlines at the concatenation point. But because "tr -cs A-Za-z '\n'" eliminates consecutive newlines, these two consecutive newlines do not appear in $f(x_1 $ ++ $ x_2)$. Therefore $f(x_1 $ ++ $ x_2) \neq g(f(x_1), f(x_2))$ and KumQuat eliminates **concat** as a potential combiner. Such counterexample inputs can be generated by input shapes whose number of words per line and number of characters per word are small.

As another example, when $f = ($uniq -c$)$ and $g = $ **concat**, a counterexample input has $x_1$ ending with a nonempty line $l$ and $x_2$ starting with the same line $l$. In this case, $f(x_1)$ ends with a line with a padded integer $n_1$ on the left and the content $l$ on the right. Meanwhile, $f(x_2)$ starts with a line with a padded integer $n_2$ on the left and the content $l$ on the right. Hence $g(f(x_1), f(x_2))$ has two consecutive lines at the concatenation point, whose contents on the right are both $l$. But because "uniq -c" merges consecutive duplicate lines, these two consecutive lines whose right sides equal do not appear in $f(x_1 $ ++ $ x_2)$. Therefore $f(x_1 $ ++ $ x_2) \neq g(f(x_1), f(x_2))$ and KumQuat eliminates **concat** as a potential combiner. Such counterexample inputs can be generated by input shapes whose percentage of distinct lines is small.

The synthesis algorithm starts with a predefined seed input shape, around which the algorithm generates a space of mutated input shapes. For each such input shape, KumQuat generates a set of input streams and feeds them to the original command. Some of these input streams may cause KumQuat to discard candidate combiners that violate the divide-and-conquer property. The sizes of the sets of discarded candidates for different input shapes induce a gradient over the input shapes. KumQuat follows this gradient to find input shapes that maximize its ability to quickly find and discard
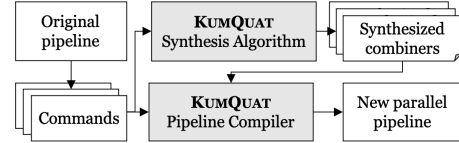


**Figure 2.** KumQuat system workflow includes splitting an original serial pipeline into commands, synthesizing combiners for each of these commands, and reassembling the commands and synthesized combiners into a new parallel pipeline.

incorrect combiners. KumQuat continues this process until it executes several gradient steps that do not discard any remaining candidate combiners. In our example this process quickly produces the correct combiners described above for the commands in our example pipeline (and typically also for the commands in our benchmark scripts, see Section 4).

**New Data-Parallel Pipeline**: KumQuat parses the original pipeline, splits it into individual commands, synthesizes combiners for these commands, and compiles them into a new data parallel pipeline (Figure 2). In our example KumQuat synthesizes combiners for all commands in the pipeline. The combiner for the command "tr -cs A-Za-z '\n'" is **rerun**. Because this command does not significantly reduce the size of the output stream (in comparison with the input stream), parallelizing the command with the **rerun** combiner reduces the overall performance. KumQuat therefore executes this command sequentially, with the input file piped directly into the command. All other commands execute in parallel. KumQuat applies the intermediate combiner elimination optimization to eliminate the combiner for the "tr A-Z a-z" command. The resulting optimized pipeline has one sequential stage and three parallel stages (one of which executes the "tr A-Z a-z" and "sort" commands with no intermediate combiner).

On a 3GB benchmark input, the serial execution time is 2089 seconds. The unoptimized parallel execution time with 16 way parallelism is 196 seconds, with a parallel speedup of 10.7×. The optimized parallel execution time is 146 seconds with a speedup of 14.4× (see Section 4).

## 3 Design

KumQuat is designed to work with commands that implement deterministic computations over input streams that are structured as units separated by delimiters. KumQuat targets commands $f$ that have a combine operator $g$ in a domain-specific combiner language that satisfies

$$f(x_1 $$ ++ $$ x_2) = g(f(x_1), f(x_2))$$

for all input streams $x_1, x_2$, where the streams terminate with newlines. KumQuat generates input streams to collect outputs from $f$, which are used to eliminate incorrect candidate combiners. Note that the combiners do not directly operate

$$
\begin{aligned}
g \in \mathsf{Combiner}_f &\ :=\ b \mid s \mid r \\
b, b_1, b_2 \in \mathsf{RecOp} &\ :=\ \mathbf{add} \mid \mathbf{concat} \mid \mathbf{first} \mid \mathbf{second} \\
&\ \mid\ \mathbf{front}\ d\ b \mid \mathbf{back}\ d\ b \mid \mathbf{fuse}\ d\ b \\
s \in \mathsf{StructOp} &\ :=\ \mathbf{stitch}\ b \mid \mathbf{stitch2}\ d\ b_1\ b_2 \mid \mathbf{offset}\ d\ b \\
r \in \mathsf{RunOp}_f &\ :=\ \mathbf{rerun}_f \mid \mathbf{merge}\ \texttt{<flags>} \\
d \in \mathsf{Delim} &\ :=\ \text{`\textbackslash n`} \mid \text{`\textbackslash t`} \mid \text{` `} \mid \text{`,`}
\end{aligned}
$$

**Figure 3.** Combiners synthesizable by KumQuat.

on the input streams generated by KumQuat, but instead operate on the command outputs. We present details of the DSL semantics, definitions, and theorems in the appendix.

### 3.1 The KumQuat Combiner DSL

To capture the space of possible combiners, KumQuat defines and uses a domain-specific language (DSL) presented in Figure 3. A combiner in KumQuat's DSL is an expression, which is a binary operation that accepts two streams $y_1$ and $y_2$ as arguments.

**Definition 3.1.** A stream is a string that ends with a newline character '\n', $\mathsf{Stream} = \{x \mathbin{+\!\!+} \text{`\textbackslash n`} \mid x \in \mathsf{String}\}$.

**Definition 3.2.** A command $f : \mathsf{Stream} \to \mathsf{Stream}$ is a function that takes a stream as input and produces a stream as output.

The DSL has three classes of operators: RecOp, StructOp, and $\mathsf{RunOp}_f$. RecOp defines recursive operators and includes numeric addition (**add**), string concatenation (**concat**), selection (**first** and **second**), and delimiter-based composite operators (**front**, **back**, **fuse**). The **front** (or **back**) operator removes a delimiter at the front (or back) of $y_1$ and $y_2$, applies a child operator, and attaches the original delimiter to the front (or back) of the combined result. The **fuse** operator applies a child operator piecewise on elements in $y_1$ and $y_2$ that are separated by a delimiter, after which piecewise results are connected back with the original delimiter.

StructOp defines operators that apply RecOp operators on structured streams (**stitch**, **stitch2**, and **offset**). These operators depend on the values at certain locations in $y_1$ and $y_2$. The **stitch** (or **stitch2**) operator compares $y_1$'s last line with $y_2$'s first line, then behaves differently conditioned on whether these two lines (or whether the second field from these two lines) equal. The **offset** operator uses the first field in the last line of $y_1$ to adjust the first field in every line of $y_2$.

$\mathsf{RunOp}_f$ defines operators that require command executions ($\mathbf{rerun}_f$ and **merge**). The $\mathbf{rerun}_f$ command reexecutes command $f$ on the concatenataion of $y_1$ and $y_2$. The **merge** command invokes a standard Unix merge command that takes two pre-sorted streams and interleaves them into a sorted merged stream. The "<flags>" parameter represents a set of known flags specific to command $f$.

Figure 6 presents several rules of the big-step execution semantics for the DSL. The transition function $\Rightarrow$ maps a DSL expression to its output value.

### 3.2 Combiner Synthesis

The synthesizer starts with an initial search space of candidate combiners in $\mathsf{Combiner}_f$. The algorithm generates a set of input streams, uses them to execute $f$, observes the outputs, and uses the observations to do two things: (1) remove implausible candidates and (2) choose an input shape for the next round of input generation. We formalize these definitions below.

**Definition 3.3.** An input pair $\langle x_1, x_2 \rangle$ consists of two strings $x_1, x_2 \in \mathsf{String}$. An output tuple $\langle y_1, y_2, y_{12} \rangle$ consists of three strings $y_1, y_2, y_{12} \in \mathsf{String}$.

**Definition 3.4.** An input stream pair $\langle x_1, x_2 \rangle$ consists of two streams $x_1, x_2 \in \mathsf{Stream}$. An observation $\langle y_1, y_2, y_{12} \rangle$ consists of three streams $y_1, y_2, y_{12} \in \mathsf{Stream}$.

**Definition 3.5.** Executing command $f$ with an input stream pair $\langle x_1, x_2 \rangle$ produces the observation $\langle f(x_1), f(x_2), f(x_1 \mathbin{+\!\!+} x_2) \rangle$. For a set of input stream pairs $X$, $f(X)$ denotes the set of observations obtained from executing $f$ with $X$, $f(X) = \{\langle f(x_1), f(x_2), f(x_1 \mathbin{+\!\!+} x_2) \rangle \mid \langle x_1, x_2 \rangle \in X\}$.

Our current KumQuat implementation allows the user to specify the initial search space with the maximum AST size in the combiner DSL.

**Definition 3.6.** The size of a combiner $g \in \mathsf{Combiner}_f$ is denoted as $|g|$ and defined as two (each combiner operates on two arguments) plus the number of times that the AST of $g$ applies a production to expand a "RecOp", "StructOp", or "$\mathsf{RunOp}_f$" symbol.

**Definition 3.7.** $G_n = \{g \in \mathsf{Combiner}_f \mid |g| \le n\}$ denotes the set of combiners that are under size $n$.

We next define legal inputs and plausible combiners.

**Definition 3.8.** For $g \in \mathsf{Combiner}_f$, $L(g)$ denotes the set of legal strings for which $g$ is defined. For example:

$$
\begin{aligned}
L(\mathbf{add}) &= [\text{`0`} - \text{`9`}]^+ \\
L(\mathbf{front}\ d\ b) &= \{d \mathbin{+\!\!+} y \mid y \in L(b)\} \\
L(\mathbf{fuse}\ d\ b) &= \{y_1 \mathbin{+\!\!+} d \mathbin{+\!\!+} y_2 \mathbin{+\!\!+} d \mathbin{+\!\!+} \ldots \mathbin{+\!\!+} d \mathbin{+\!\!+} y_k \\
&\quad \mid y_1 \ne \mathsf{nil},\ y_k \ne \mathsf{nil},\ \text{and}\ y_i \in L(b)\ \text{and}\ d \notin y_i \\
&\quad \text{for all}\ i = 1, \ldots, k,\ \text{where}\ k \ge 2\}
\end{aligned}
$$

For any $y_1, y_2 \in L(g)$, the evaluation $g\ y_1\ y_2 \Longrightarrow_e v$ succeeds for some $v \in \mathsf{String}$.

**Definition 3.9.** A combiner $g \in \mathsf{Combiner}_f$ is plausible for output tuples $Y$, denoted as $P(g, Y)$, if $y_1, y_2 \in L(g)$ and $g\ y_1\ y_2 \Longrightarrow_e y_{12}$ for all $\langle y_1, y_2, y_{12} \rangle \in Y$.

**Definition 3.10.** A combiner $g \in \mathsf{Combiner}_f$ is correct with respect to input pairs $X$ if $P(g, f(X))$. $G_n^{f,X} = \{g \in G_n \mid P(g, f(X))\}$ denotes the combiners for command $f$ that are under size $n$ and are correct with respect to input pairs $X$.

$$\frac{i_1 = \text{strToInt } y_1 \qquad i_2 = \text{strToInt } y_2}{\textbf{add } y_1 \ y_2 \Longrightarrow_e \text{intToStr } (i_1 + i_2)}$$

$$\frac{h_1, t_1 = \text{splitFirst } d \ y_1 \qquad h_2, t_2 = \text{splitFirst } d \ y_2 \qquad t_1 \neq \text{nil} \qquad t_2 \neq \text{nil}}{d \in t_1 \qquad d \in t_2 \qquad b \ h_1 \ h_2 \Longrightarrow_e v \qquad (\textbf{fuse } d \ b) \ t_1 \ t_2 \Longrightarrow_e v'}$$

$$\overline{(\textbf{fuse } d \ b) \ y_1 \ y_2 \Longrightarrow_e v +\!\!+ \ d +\!\!+ \ v'}$$

$$\frac{}{\textbf{concat } y_1 \ y_2 \Longrightarrow_e y_1 +\!\!+ y_2}$$

$$\frac{y'_1, l_1 = \text{splitLastLine } y_1 \qquad l_2, y'_2 = \text{splitFirstLine } y_2 \qquad l_1 = l_2 \qquad b \ l_1 \ l_2 \Longrightarrow_e v}{(\textbf{stitch } b) \ y_1 \ y_2 \Longrightarrow_e y'_1 +\!\!+ \ `\backslash n` +\!\!+ v +\!\!+ \ `\backslash n` +\!\!+ y'_2}$$

$$\frac{}{\textbf{first } y_1 \ y_2 \Longrightarrow_e y_1}$$

$$\frac{b \ (\text{delFront } d \ y_1) \ (\text{delFront } d \ y_2) \Longrightarrow_e v}{(\textbf{front } d \ b) \ y_1 \ y_2 \Longrightarrow_e d +\!\!+ v}$$

$$\frac{v = (\text{unixMerge <flags>}) \ y_1 \ y_2}{(\textbf{merge <flags>}) \ y_1 \ y_2 \Longrightarrow_e v} \qquad \frac{}{\textbf{rerun}_f \ y_1 \ y_2 \Longrightarrow_e f \ (y_1 +\!\!+ y_2)}$$

$b, b_1, b_2 \in \text{RecOp} \qquad d \in \text{Delim} \qquad y_1, y_2, y'_1, y'_2, v, v', v_1, v_2, h, h_1, h_2, t, t_1, t_2, l_1, l_2 \in \text{String} \qquad i_1, i_2 \in \text{Int}$

**Figure 4.** Semantics of the combiner DSL (selected rules).

**Data:** Command $f$, max combiner size $n$
**Result:** Synthesized plausible combiners
$C_0 \leftarrow AllCandidates(n)$
**for** $r = 1, 2, \ldots$ **do**
  $I_r \leftarrow GetEffectiveInputs(f, C_{r-1}, RandomShape())$
  $C_r \leftarrow FilterCandidates(f, C_{r-1}, I_r)$
  **if** $C_r = \emptyset$ **then**
    | **return** nil
  **end**
  **if** *not* $MakingProgress([C_0, \ldots, C_r])$ **then**
    | **return** $C_r$
  **end**
**end**

**Algorithm 1:** Procedure *Synthesize*, which implements KumQuat's core synthesis algorithm. The procedure takes a command and synthesizes a combiner for the command that is correct with respect to a range of generated input pairs.

**Data:** Command $f$, candidate combiners $C$, input shape $s_0$
**Result:** Input stream pairs, generated from mutating $s_0$, for eliminating incorrect candidates in $C$
$I \leftarrow$ Empty set
**for** $m = 1, \ldots, M$ **do**
  **for** $j = 1, \ldots, 12$ **do**
    $s^j_{m-1} \leftarrow MutateShape(s_{m-1}, j)$
    $I^j_{m-1} \leftarrow GetInputStreamPairs(s^j_{m-1})$
    Add $I^j_{m-1}$ to $I$
  **end**
  $j' \leftarrow IndexBestMutation(C, I^1_{m-1}, \ldots, I^{12}_{m-1})$
  $s_m \leftarrow s^{j'}_{m-1}$
**end**
**return** $I$

**Algorithm 2:** Procedure *GetEffectiveInputs*, which mutates input shapes to generate input stream pairs.

**Combiner Synthesis**: Algorithm 1 presents KumQuat's combiner synthesis algorithm, which takes a black-box command $f$ and an integer $n$. It starts by preparing a set $C_0$ of the initial search space. The algorithm then performs multiple rounds of filtering on the candidates.

Variable $C_r$ holds the set of combiners that are correct with respect to all seen input pairs, *i.e.*, the set $G^{f,I}_n$ where $I = \bigcup^r_{r'=1} I_{r'}$. Algorithm 1 terminates if either (1) no candidate combiners remain, in which case it returns nil and reports an error, or (2) no progress is made in a number of rounds, in which case it returns the set of plausible combiners.

**Input Generation**: We next present how the inputs are generated for Algorithm 1 to filter candidates. A key goal of input generation is to generate a variety of input streams that exercise a wide range of the functionality of the command $f$. The KumQuat input generation algorithm is driven by mutations to an input shape, from which KumQuat generates random inputs. The mutations are chosen by how effectively their resulting inputs eliminate incorrect candidate combiners.

**Definition 3.11.** An *input shape* $s = \langle s_L, s_W, s_C \rangle \in Shape$ specifies the configurations for three *dimensions* of an input: the lines in each input as separated by newline characters ($s_L \in Config$), the words in each line as separated by spaces ($s_W \in Config$), and the characters in each word ($s_C \in Config$). The configuration for each dimension is of the form $\langle l, u, d \rangle$ and specifies three bounds: the minimum element count ($l \in Int$), the maximum element count ($u \in Int$), and the percentage of distinct elements ($d \in Percent$) on that dimension.

**Definition 3.12.** A stream $x$ *satisfies* an input shape $s \in Shape$, denoted as $x \sim s$, if $x$ conforms to the bounds specified in $s$. An input stream pair $\langle x_1, x_2 \rangle$ satisfies an input shape $s \in Shape$, denoted as $\langle x_1, x_2 \rangle \sim s$, if $(x_1 +\!\!+ x_2) \sim s$.

Algorithm 2 presents KumQuat's input generation algorithm. The procedure takes a black-box command $f$, a set of candidate combiners $C$, and an initial input shape $s_0$; it mutates the input shape iteratively, generating input streams along the way.

The iterative mutation process is inspired by gradient descent. The $m$-th iteration ($m = 1, \ldots, M$) mutates the input shape $s_{m-1}$ using one of twelve potential mutations. These

potential mutations are along three dimensions (lines, words, and characters) and four directions (more/fewer elements, more/less varied). Procedure *MutateShape* takes an initial input shape and a mutation index, then returns a new input shape mutated as specified. For the $j$-th potential mutation ($j = 1, \ldots, 12$), Algorithm 2 uses the mutated input shape $s_{m-1}^j$ to generate a set of input stream pairs $I_{m-1}^j$. In other words, the variable $I_{m-1}^j$ satisfies $\langle i_1, i_2 \rangle \sim s_{m-1}^j$ for all $\langle i_1, i_2 \rangle \in I_{m-1}^j$.

Algorithm 2 then evaluates the effectiveness of all of the input shape mutations. It returns the index, $j'$, of the most effective set. The $j'$-th mutation then produces the input shape for the next iteration, $s_m = s_{m-1}^{j'}$. The procedure repeats these operations for $M$ iterations. Finally the procedure returns the set of all observed input pairs, $I = \bigcup_{j,m} I_{m-1}^j$.

**Preprocessing**: The current KumQuat implementation preprocesses command scripts to obtain a set of literals for generating inputs and input shapes. For example, "grep 'light.light'" does not produce any outputs unless the input stream contains lines that match the regular expression 'light.light'. KumQuat extracts this regular expression and generates a dictionary of strings that match. It then uses this dictionary as elements for generating input streams based on input shapes. The command "sed 100q" copies input to output when the input stream contains at most 100 lines. When the input stream contains more lines, the command removes the trailing lines. KumQuat obtains the number 100 as a literal and uses it to generate initial input shapes where one dimension is close to this number. KumQuat then mutates this initial input shape to obtain a range of different input streams that exercise different behavior in the original command.

KumQuat also checks whether the original command can process three test input streams without errors: a list of unsorted English words separated by newlines, the same list of words but sorted, and a list of legal file names separated by newlines. Most of our benchmark commands can process all three test input streams without errors. Benchmark commands that use "comm" print an error with the first test input stream but succeed with the second. Based on this outcome, KumQuat generates only sorted input streams for these commands during combiner synthesis. Benchmark commands that use "xargs" print an error with the first two test input streams but succeed with the third. Based on this outcome, KumQuat configures a dictionary of legal file names and uses the dictionary to generate input streams for these commands.

**Multiple Plausible Combiners**: Recall that Algorithm 1 returns the set of plausible combiners (Definition 3.9 and Definition 3.10). Let $G$ be the set of returned plausible combiners.

If $G = \emptyset$, the synthesizer reports an error. If $G$ contains exactly one combiner, the synthesizer returns the combiner directly.

If $G$ contains more than one combiner, the synthesizer builds a composite combiner using the following subset of $G$. If $G \cap \text{RecOp} \neq \emptyset$, the synthesizer uses the set $G \cap \text{RecOp}$. Otherwise, if $G \cap \text{StructOp} \neq \emptyset$, the synthesizer uses the set $G \cap \text{StructOp}$. Otherwise, $G \cap \text{RunOp}_f \neq \emptyset$ and the synthesizer uses the set $G \cap \text{RunOp}_f$. Let this nonempty subset be $\{g_1, g_2, \ldots, g_m\}$ ($m \geq 1$). KumQuat uses it to construct a composite combiner as follows. For any $y_1, y_2$, if they belong to the domain of $g_1$, then return $g_1(y_1, y_2)$. Else if they belong to the domain of $g_2$, return $g_2(y_1, y_2)$. . . . Otherwise, return $g_m(y_1, y_2)$. We show in Section 3.3 that, if the correct combiner for $f$ is among a certain set, then the order in which these combiners are composed together does not matter—the resulting composite combiner is semantically equivalent regardless of the order. Alternatively, if the domain of one of these plausible combiners is the superset of any other plausible combiner's domain, then it suffices to return only the combiner with the largest domain.

## 3.3 Conditions for Synthesizing Correct Combiners

We present theorems that characterize when the combiner synthesis algorithm will identify a correct combiner for a given set of parallel output streams.

Broadly speaking, when the combiner involves numerical addition (**add**), the corresponding stream fragments on which the numerical addition applies are required to be nonzero in some observations. When the combiner involves string concatenation (**concat**), the corresponding fragments on which the string concatenation applies are required to be nonempty in some observations. When the combiner involves selection (**first** and **second**), the corresponding stream fragments on which the selection applies are required to contain non-delimiter and non-zero characters in some observations.

Combiners that process formatted streams may nest these three classes of basic operators inside more complex operators (**front**, **back**, **fuse**, **stitch**, **stitch2**, **offset**). For these combiners, their requirements for sufficient observations include a specification of the formatting as well as a specification of the deformatted fragments.

**Definition 3.13.** For $g_1, g_2 \in \text{Combiner}_f$, $g_1$ and $g_2$ are equivalent by intersection, denoted as $g_1 \equiv_\cap g_2$, if for all $y_1, y_2 \in L(g_1) \cap L(g_2)$, $g_1 \, y_1 \, y_2 \Longrightarrow_e v$ and $g_2 \, y_1 \, y_2 \Longrightarrow_e v$ for some $v$.

**Definition 3.14.** A combiner $g \in \text{Combiner}_f$ is correct for command $f$ if $P(g, f(X))$ holds for all input stream pairs $X$.

**Definition 3.15.** We define two sets of representative combiners for command $f$, $G_{\text{rec}} = \{g_a, g_c, g_f, g_s, g_{ba}, g_{fa}, g_{bfa}, g_{fbfa}, g_{fc}\} \subset \text{RecOp}$ and $G_{\text{struct}} = \{g_{sf}, g_{saf}, g_{oa}\} \subset \text{StructOp}$, whose

elements include: $g_{\mathrm{c}}$ = **concat**, $g_{\mathrm{ba}}$ = (**back** $d$ **add**), $g_{\mathrm{sf}}$ = (**stitch first**), and $g_{\mathrm{saf}}$ = (**stitch2** $d$ **add first**).

**Definition 3.16.** For combiner $g \in G_{\mathrm{rec}} \cup G_{\mathrm{struct}}$ and any set of output tuples $Y$, $E(g, Y)$ denotes a conservative predicate that is true only if $Y$ is sufficient for eliminating incorrect candidates when the correct combiner is $g$.

**Definition 3.17.** For any set of output tuples $Y$, $E_{\mathrm{rec}}(Y)$ denotes a conservative predicate that is true only if $Y$ is sufficient for eliminating incorrect candidates when the correct combiner $g \in G_{\mathrm{rec}}$.

**Definition 3.18.** For any set of output tuples $Y$, $T(Y)$ denotes a predicate that is true only if $Y$ is interpretable as a table.

**Definition 3.19.** For any set of output tuples $Y$, $E_{\mathrm{struct}}(Y)$ denotes a conservative predicate that is true only if $Y$ is sufficient for eliminating incorrect candidates when the correct combiner $g \in G_{\mathrm{struct}}$.

**Definition 3.20.** For a command $f$, integer $k$, and set of output tuples $Y$, the set of plausible combiners $P_k(Y) = \{g \in \mathrm{Combiner}_f \mid |g| \le k$ and $P(g, Y)\}$.

**Theorem 1.** For any combiner $g \in G_{\mathrm{rec}}$, set of output tuples $Y$ such that $P(g, Y)$ and $E(g, Y)$, and $g' \in \mathrm{RecOp}$, we have $P(g', Y)$ implies $g' \equiv_\cap g$.

**Theorem 2.** For any command $f$, set of input streams $X$, combiner $g \in G_{\mathrm{rec}}$, combiner $g' \in \mathrm{RecOp}$, and integer $k$, if the following conditions hold:

- $E_{\mathrm{rec}}(f(X))$,
- $g$ is correct for $f$,
- $y_1, y_2 \in L(g')$ for all $\langle y_1, y_2, y_{12} \rangle \in f(X)$, and
- $k \ge |g'|$,

then $g' \in P_k(f(X)) \cap \mathrm{RecOp}$ if and only if $g' \equiv_\cap g$.

**Theorem 3.** For any combiner $g \in G_{\mathrm{struct}}$, set of output tuples $Y$ such that $P(g, Y)$ and $E(g, Y)$, and $g' \in \mathrm{StructOp}$, we have $P(g', Y)$ implies $g' \equiv_\cap g$.

**Theorem 4.** For any command $f$, set of input streams $X$, combiner $g \in G_{\mathrm{struct}}$, combiner $g' \in \mathrm{StructOp}$, and integer $k$, if the following conditions hold:

- $E_{\mathrm{struct}}(f(X))$,
- $g$ is correct for $f$,
- $y_1, y_2 \in L(g')$ for all $\langle y_1, y_2, y_{12} \rangle \in f(X)$, and
- $k \ge |g'|$,

then $g' \in P_k(f(X)) \cap \mathrm{StructOp}$ if and only if $g' \equiv_\cap g$.

### 3.4 Input Generation and Correct Combiners

Our target commands often consist of two components: unit-based computation and delimiter-based formatting. Unit-based computation often determines how a combiner applies the **add**, **concat**, **first**, and **second** operators to certain fragments of the output substreams. Delimiter-based formatting often determines how a combiner uses the **front**, **back**, **fuse**,

**stitch**, **stitch2**, and **offset** operators. We focus on three broad classes of unit-based computation that appear in our benchmark commands.

**Counting Lines, Words, or Characters**: Many Uɴɪx commands output formatted counts of certain lines, words, or characters. Benchmark commands that implement this pattern include "wc -l," "grep -c [regex]," and "uniq -c."

For each of these wc and grep commands, a correct combiner is (**back** $d$ **add**). By Theorem 1 and Theorem 2, as long as the command outputs satisfy the requirements for (**back** $d$ **add**), any synthesized plausible combiner in RecOp will be equivalent to (**back** $d$ **add**) when processing streams that belong to the combiner's domain. For the uniq command, a correct combiner is (**stitch2** ' ' **add first**). By Theorem 3 and Theorem 4, as long as the command outputs collected by KᴜᴍQᴜᴀᴛ satisfy the requirements for (**stitch2** ' ' **add first**), any synthesized plausible combiner in StructOp will be equivalent to (**stitch2** ' ' **add first**) when processing streams that belong to the combiner's domain.

Both of these correct combiners use **add** nested inside other operators. Because these other operators process formatting, the remaining deformatted fragments in the output substreams are therefore processed by **add**. Here we focus on identifying the correct **add** operator for processing these deformatted fragments. The requirement is (conceptually) observing nonzero values in these fragments.

KᴜᴍQᴜᴀᴛ generates input streams with various numbers of lines. Many of these lines cause the wc counter(s) to be nonzero. For "grep -c [regex]," the KᴜᴍQᴜᴀᴛ preprocessing extracts literals that it uses to generate input streams that contain matching values that cause the grep counter to be nonzero. The resulting output streams therefore contain nonzero characters even after removing the formatting, which satisfy the requirements for identifying **add** correctly as a building block of the final synthesized combiners.

For "uniq -c" the combiner contains a conditional in the **stitch2** operator that applies the **add** operator only when the right-hand content in the last line of $y_1$ equals the right-hand content in the first line of $y_2$. These contents correspond to the last line of $x_1$ and the first line of $x_2$. KᴜᴍQᴜᴀᴛ generates input streams $x_1, x_2$ with varying percentages of distinct lines, some of which enable the command to produce output streams $y_1, y_2$ that contain deformatted fragments that are processed by the **add** operator. Because these deormatted fragments are always nonzero for this uniq command, they satisfy the requirements for identifying **add** correctly as a building block of the final synthesized combiner.

**Mapping Input Lines to Disjoint Output Lines**: Many Uɴɪx commands apply a function to map each input line to a sequence of output lines. Benchmark commands that implement this pattern include: "tr '[a-z]' 'P'," "tr -c '[A-Z]' '\n'," "sed s/\$/'0s'/," "cut -c 1-4," "cut -d

',' -f 3,1," "awk "length >= 16"," "grep 'light.\*light',"
"grep -v '^0$'," "xargs cat," and "xargs file."

For each of these commands, a correct combiner is **concat**. By Theorem 1 and Theorem 2, as long as the command outputs collected by KumQuat satisfy the requirements for **concat**, any synthesized plausible combiner in RecOp will be equivalent to **concat** when processing streams that belong to the combiner's domain. In this case, the requirement is (conceptually) observing nonempty output streams.

For the tr, sed, and cut commands, KumQuat generates input streams with various numbers of lines, many of which cause these commands to produce nonempty outputs. For the awk, grep, and xargs commands, KumQuat uses pre-processing to determine that it will generate input streams based on certain literals or file names (Definition 3.2). The resulting input streams therefore enable these commands to produce nonempty outputs, which satisfy the requirements for synthesizing **concat** correctly.

**Selecting Elements**: Some Unix commands select certain elements from a list, output selected elements, and discard others. Benchmark commands include: "uniq" and "uniq -c." (**stitch first**) is a correct combiner for "uniq"; (**stitch2** ' ' **add first**) is a correct combiner for "uniq -c".

These combiners use **first** or **second** nested inside deformatting operators so that correct fragments are processed by **first** or **second**. Here we focus on how to identify the correct **first** or **second** operator. The requirement is (conceptually) observing non-delimiter and non-zero characters in these fragments and observing such fragments for the two operands to differ.

These uniq commands select one line out of any two adjacent lines that equal. The combiners contain conditional statements in the **stitch** and **stitch2** operators that apply **first** or **second** only when certain contents are equal. The combiner for "uniq" applies the **first** operator when the last line in $y_1$ equals the first line in $y_2$. The combiner for "uniq -c" applies the **first** operator when the right-hand content in the last line of $y_1$ equals the right-hand content in the first line of $y_2$. KumQuat generates input streams $x_1, x_2$ with varying percentages of distinct lines, some of which enable these commands to produce output streams $y_1, y_2$ that contain deformatted fragments that are processed by the **first** operator. Also, because KumQuat generates input streams with varying numbers of words per line, some of the generated lines contain nonzero and nondelimiter characters which enable the command to produce such characters in the deformatted fragments as well. These deformatted fragments therefore satisfy the requirements for identifying **first** and **second** as a building block of the final synthesized combiner.

### 3.5 Pipeline Optimization

**Eliminating Intermediate Combiners**: KumQuat eliminates unnecessary intermediate combiners as follows.



**(a)** Serial pipeline that consists of two commands $f_1, f_2$



**(b)** Unoptimized parallel pipeline that executes a combiner after each parallel command



**(c)** Optimized parallel pipeline that executes multiple commands in parallel after eliminating intermediate combiners

**Figure 5.** KumQuat reassembles a new parallel pipeline by splitting the input stream, running parallel copies of the original commands on the input substreams, and combining the output substreams.

**Theorem 5.** For any commands $f_1, f_2$, input streams $x_1, x_2 \in$ Stream, correct combiner $g_1$ for $f_1$, correct combiner $g_2$ for $f_2$ if $g_1 = $ **concat** and $f_1(x_1), f_1(x_2) \in$ Stream then

$$g_2(f_2(y_1), f_2(y_2)) = g_2(f_2(f_1(x_1)), f_2(f_1(x_2)))$$

for any $y_1, y_2 \in$ Stream such that $f_1(x_1 ++ x_2) = y_1 ++ y_2$.

Figure 5 shows the effect of this optimization—rather than combining after every pipeline stage, the optimized parallel pipeline combines only once. In general, if the combiner for a stage (command $f_1$) concatenates the parallel output substreams, then we can eliminate the combiner $g_1$ for this stage and feed the output substreams directly to the next parallel stage $f_2$ as input substreams. The final combined output stream applies only the combiner $g_2$ for the second stage (Figure 5c). By Theorem 5, this output is identical to the unoptimized output if $g_1$ were in place (Figure 5b).

A prerequisite for this optimization is that the command $f_1$ must produce output streams that terminate with newlines. One of our benchmark commands violate this precondition: the command "tr -d '\n'" removes all newline characters. Hence the optimization in Theorem 5 does not apply to this command (KumQuat still parallelizes this command with the **concat** combiner).

**Combining Multiple Substreams**: Although KumQuat synthesizes combiners that process two streams, the commands may be executed with $k$ way parallelism that produces $k$ output substreams ($k > 2$). KumQuat generalizes the following combiners to apply to all $k$ substreams at the same time. The **merge** <flags> combiner is implemented in KumQuat as an invocation of a Unix script "sort -m <flags> $*" which merges multiple sorted streams at the same time. The **concat** combiner is implemented as the script "cat $*" which concatenates multiple streams. The **rerun**

combiner can also be implemented by concatenating all sub-streams at the same time and rerunning the original command only once. For other combiners, the current KumQuat implementation applies the combiner on two substreams repeatedly until only one substream remains.

## 4 Experimental Results

We evaluate KumQuat on the following benchmarks:

- **Mass-transit analytics during COVID-19:** This benchmark set contains 4 scripts that were used to analyze real telemetry data from bus schedules during the COVID-19 response in a large European city [25]. The pipelines compute several average statistics on the transit system per day—such as daily serving hours and daily number of vehicles. Each script has 1 pipeline. Each pipeline has between 7 and 8 stages.[3] These scripts operate on a fixed 3.4GB dataset that contains mass-transport data collected over a single year.

- **Natural language processing:** This benchmark set contains 22 scripts from Kenneth's Unix-for-Poets [5], updated in 2016 by a Stanford linguistic class [15]. These scripts calculate natural-language processing metrics such as n-grams, morphs, counts, and frequencies. Each script has between 1 and 3 pipelines. Each pipeline has between 2 and 8 stages. These scripts are applied to 1823 books that total 927MB from Project Gutenberg [11].

- **Classic Unix One-liners:** This benchmark set contains 10 pipelines written by Unix experts: a few pipelines are from Unix legends [1, 2, 16], one from a book on Unix scripting [24], and a few are from top Stackoverflow answers [13]. Each script has between 1 and 2 pipelines, except for a script that has only one command. Each pipeline has between 2 and 8 stages. Inputs are script-specific and average 1.6GB per benchmark.

- **Unix50 from Bell Labs:** This benchmark set contains 34 pipelines solving the Unix 50 game [14], designed to highlight Unix's modular philosophy [16], found on GitHub [3] Each script has 1 pipeline, except for a script that has only one command. Each pipeline has between 2 and 10 stages. Inputs are script-specific and average 1.1GB per benchmark.

**Experimental Setup**: To evaluate the pipeline performance, we implemented an infrastructure that can execute each stage in a pipeline to completion before starting to execute the next stage. The infrastructure configures any stage that invokes the Unix `sort` utility to be serial (using the option "`--parallel=1`"). Each stage's output is redirected to a file, which is read by the next stage as input. This infrastructure

provides a parameter for specifying the amount of parallelism for each parallelizable stage.

We performed experiments on a server with 0.5TB of memory and $80 \times 2.27$GHz Intel(R) Xeon(R) E7-8860, Debian GNU/Linux 9, GNU Coreutils 8.26-3, and Python 3.8.2. We note that our benchmarks never come close to exhausting the server's available memory.

**Performance Results**: The 70 benchmark scripts have a total of 477 commands and 427 pipeline stages.[4] Table 1 presents the performance results for our automatically parallelized pipelines for the two longest-running scripts in each benchmark. In general, shorter scripts have smaller parallel speedup. We present full results in the appendix.

The first two columns present the benchmark and script names. The next column (**Parallelized**) presents the number of stages automatically parallelized by KumQuat, $k$, and the number of stages in the original pipeline, $n$, as a pair "$k/n$" for each pipeline in the parentheses. Here we also report the single commands that are not in pipelines, as pairs "$k/1$". The pair before the parentheses presents the sum over all pipelines in the script. The next column (**Eliminated**) presents the number of parallelized stages whose combiners are eliminated by KumQuat during optimization. Again, the numbers in the parentheses correspond to pipelines in the script. The number before the parentheses presents the sum over all pipelines. Among all benchmark scripts, KumQuat parallelizes 325 of the the 427 stages (76.1%) with synthesized combiners. The optimization eliminates 144 of these combiners (44.3%). These results highlight the ability of KumQuat to effectively extract the parallelism implicitly present in the benchmark pipelines.

The next column ($T_{\mathbf{orig}}$) presents the execution time for the original benchmark script, which exploit the default Unix pipelined parallelism and deploy the default Unix sort, which exploits 8 way parallelism. The next two columns ($u_1$ and $u_{16}$) present the execution time for the *unoptimized* pipeline with 1 and 16 way parallelism for each data parallel command. Since these generated pipelines always wait for each stage to terminate before starting the next stage, $u_1$ is the serial execution time. The last column ($T_{16}$) presents the execution time for the *optimized* pipeline with 16 way command parallelism. Among the benchmark scripts whose serial execution time is at least 3 minutes, the unoptimized parallel speedup ranges between 3.5× and 14.9×, with a median speedup of 8.5×. Among these benchmark scripts, the optimized parallel speedup ranges between 3.8× and 26.9×, with a median speedup of 11.3× (we attribute the superlinear speedup to pipelined parallelism exploited across consecutive parallelized commands with no intermediate combiner).

**Synthesis Results**: We summarize the synthesis results below and present full results in the appendix. The benchmarks

---

[3]We count pipelines as groups of two or more commands connected by Unix pipes. We count pipeline stages as commands in the pipeline excluding initial "`cat`" commands that read input files.

[4]See footnote 3.

**Table 1.** Performance results for the two longest-running scripts from each benchmark

| Benchmark | Script Name | Parallelized | Eliminated | $T_{\mathbf{orig}}$ | $u_1$ | $u_{16}$ | $T_{16}$ |
|---|---|---|---|---|---|---|---|
| analytics-mts | 2.sh (vehicle days on road) | 8/8 (8/8) | 3 (3) | 335 s (1.1×) | 379 s | 41 s (9.3×) | 28 s (13.5×) |
| analytics-mts | 3.sh (vehicle hours on road) | 8/8 (8/8) | 3 (3) | 408 s (1.0×) | 427 s | 51 s (8.4×) | 38 s (11.3×) |
| oneliners | set-diff.sh | 5/8 (0/1, 3/3, 2/2, 0/1, 0/1) | 3 (0, 2, 1, 0, 0) | 879 s (1.5×) | 1308 s | 144 s (9.1×) | 128 s (10.2×) |
| oneliners | wf.sh | 4/5 (4/5) | 1 (1) | 1155 s (1.8×) | 2089 s | 196 s (10.7×) | 145 s (14.4×) |
| poets | 4_3b.sh (count_trigrams) | 4/9 (2/4, 0/1, 0/1, 2/3) | 1 (1, 0, 0, 0) | 862 s (1.2×) | 1049 s | 275 s (3.8×) | 279 s (3.8×) |
| poets | 8.2_2.sh (bigrams_appear_twice) | 4/9 (2/4, 0/1, 2/3, 0/1) | 1 (1, 0, 0, 0) | 645 s (1.4×) | 921 s | 177 s (5.2×) | 91 s (10.2×) |
| unix50 | 21.sh (8.4: longest words w/o hyphens) | 3/3 (3/3) | 1 (1) | 428 s (1.7×) | 733 s | 64 s (11.4×) | 49 s (14.9×) |
| unix50 | 23.sh (9.1: extract word PORT) | 6/6 (6/6) | 4 (4) | 111 s (1.8×) | 202 s | 23 s (8.8×) | 10 s (19.8×) |

contain 133 unique command/flag combinations (we refer to them as "commands" below). Among these commands, 121 are data-processing commands that read an input stream.[5] KumQuat synthesizes a combiner for 113 of the 121 unique commands, with no combiner synthesized for the remaining 8 commands. The 8 unsupported commands include 7 commands for which no correct combiner exists and 1 command that requires a specific field of the input file to equal "2".

The synthesis times vary between 39 seconds and 331 seconds with a median of 60 seconds. The most common synthesized plausible combiners, including their equivalents, are: **concat** (synthesized 81 times), **rerun** (30 times), **merge**($\star$) (16 times), and (**back** '\n' **add**) (12 times). Other synthesized combiners involve operators **first**, **second**, **fuse**, **stitch**, and **stitch2**. For each benchmark command, the synthesized plausible combiners are all equivalent when operating on the command's outputs. KumQuat uses the synthesized combiners to parallelize the benchmark scripts. The generated parallel pipelines all produce correct outputs (same outputs as the original scripts).

## 5 Related Work

We discuss related work in parallel execution of shell commands and scripts, synthesis of divide-and-conquer computations, program synthesis driven by provided input/output examples, and synthesis of Unix commands.

**POSH and PaSh**: The POSH and PaSh systems parallelize and distribute Unix shell scripts [18, 26]. Both systems require combiners and both systems work with manually coded combiners. KumQuat eliminates the need for manually coded combiners, enabling such systems to immediately work with new commands (or new combinations of command flags) that require new combiners without the need to manually develop new combiners.

**Synthesis of MapReduce Programs from Examples**: [21] present a technique for automatically synthesizing complete

MapReduce programs given a partial specification in the form of a set of input/output examples. KumQuat, in contrast, supports (but does not synthesize) commands with much more sophisticated semantics than the synthesized map computations in [21]. By working with existing shell commands, KumQuat also eliminates the need for the user to provide input/output examples.

Preserving the order in which components appear in output streams is required to correctly implement the streaming semantics of Unix pipelines. KumQuat preserves this required order by incorporating metadata into the combiners and the parallelization. [21], in contrast, does not support ordered streams — it targets computations that do not have ordering constraints and can produce output components in any order.

**Automatic Parallelization of Divide and Conquer Computations**: Some research in this area uses program analysis, typically over loops that access dense arrays or matrices, to generate parallel divide and conquer computations [10, 19]. Other research works with a complete characterization of the semantics of the original sequential computation [7, 8]. KumQuat, in contrast, synthesizes combiners for black-box streaming computations. KumQuat can therefore successfully target much more complex computations implemented in arbitrary programming languages. A trade-off is that the correctness of the KumQuat combiner synthesis algorithm relies on assumptions about the computation that the black-box components implement.

**Unix Synthesis**: Prior work on synthesis for Unix shell commands and pipelines [4, 6] is guided by examples or natural-language specifications. Instead of automatically generating parallel or distributed versions of an existing command or pipeline, the goal is to synthesize the sequential command itself from examples or natural language specs.

**Commands and Shells**: There is a series of systems that aid developers in running commands or script fragments in a parallel or distributed fashion. These range from simple Unix utilities [20, 23, 27] to parallel/distributed shells [17,

---

[5]The remaining 12 unique commands include 2 function calls, 3 commands that do not process data streams (ls, mkfifo, and rm), and 7 commands that process multiple input streams.

22] to data-parallel frameworks that incorporate Unix commands [9, 12] These tools require developers to modify programs to make use of the tools' APIs. KumQuat, in contrast, aims to provide an automated solution that works directly on sequential scripts.

## 6 Conclusion

KumQuat synthesizes combiners that enable the exploitation of data parallelism in Unix commands and pipelines. Our experimental results show that the KumQuat input generation and combiner synthesis algorithms effectively identify correct combiners for our benchmark scripts and that these combiners enable the effective parallelization of these scripts.

## Acknowledgments

## References

[1] Jon Bentley. 1985. Programming Pearls: A Spelling Checker. *Commun. ACM* 28, 5 (May 1985), 456–462. https://doi.org/10.1145/3532.315102

[2] Jon Bentley, Don Knuth, and Doug McIlroy. 1986. Programming Pearls: A Literate Program. *Commun. ACM* 29, 6 (June 1986), 471–483. https://doi.org/10.1145/5948.315654

[3] Pawan Bhandari. 2020. *Solutions to unixgame.io.* https://git.io/Jf2dn Accessed: 2020-04-14.

[4] Sanjay Bhansali and Mehdi T Harandi. 1993. Synthesis of UNIX programs using derivational analogy. *Machine Learning* 10, 1 (1993), 7–55.

[5] Kenneth Ward Church. 1994. Unix™for poets. *Notes of a course from the European Summer School on Language and Speech Communication, Corpus Based Methods* (1994).

[6] Anthony Cozzie, Murph Finnicum, and Samuel T King. 2011. Macho: Programming with Man Pages. In *13th Workshop on Hot Topics in Operating Systems.* USENIX Association, Napa, CA, United States. http://www.usenix.org/events/hotos11/tech/final_files/Cozzie.pdf

[7] Azadeh Farzan and Victor Nicolet. 2017. Synthesis of Divide and Conquer Parallelism for Loops. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Barcelona, Spain) *(PLDI 2017).* Association for Computing Machinery, New York, NY, USA, 540–555. https://doi.org/10.1145/3062341.3062355

[8] Azadeh Farzan and Victor Nicolet. 2019. Modular Divide-and-Conquer Parallelization of Nested Loops. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Phoenix, AZ, USA) *(PLDI 2019).* Association for Computing Machinery, New York, NY, USA, 610–624. https://doi.org/10.1145/3314221.3314612

[9] Apache Software Foundation. 2020. Hadoop Streaming. https://hadoop.apache.org/docs/r1.2.1/streaming.html

[10] Manish Gupta, Sayak Mukhopadhyay, and Navin Sinha. [n.d.]. Automatic Parallelization of Recursive Procedures. In *Proceedings of the 1999 International Conference on Parallel Architectures and Compilation Techniques, Newport Beach, California, USA, October 12-16, 1999.* 139–148.

[11] Michael Hart. 1971. *Project Gutenberg.* https://www.gutenberg.org/

[12] Michael Isard, Mihai Budiu, Yuan Yu, Andrew Birrell, and Dennis Fetterly. 2007. Dryad: distributed data-parallel programs from sequential building blocks. In *Proceedings of the 2nd ACM SIGOPS/EuroSys European Conference on Computer Systems 2007.* 59–72.

[13] Dan Jurafsky. 2017. Unix for Poets. https://web.stanford.edu/class/cs124/lec/124-2018-UnixForPoets.pdf

[14] Nokia Bell Labs. 2019. *The Unix Game—Solve puzzles using Unix pipes.* https://unixgame.io/unix50 Accessed: 2020-03-05.

[15] Christopher Manning. 2016. *Unix for Poets (in 2016).* https://web.stanford.edu/class/archive/linguist/linguist278/linguist278.1172/notes/278-UnixForPoets.pdf Accessed: 2021-03-11.

[16] Malcolm D McIlroy, Elliot N Pinson, and Berkley A Tague. 1978. UNIX Time-Sharing System: Foreword. *Bell System Technical Journal* 57, 6 (1978), 1899–1904.

[17] Rob Pike, Dave Presotto, Ken Thompson, Howard Trickey, et al. 1990. Plan 9 from Bell Labs. In *Proceedings of the summer 1990 UKUUG Conference.* 1–9. http://css.csail.mit.edu/6.824/2014/papers/plan9.pdf

[18] Deepti Raghavan, Sadjad Fouladi, Philip Levis, and Matei Zaharia. 2020. POSH: A Data-Aware Shell. In *2020 USENIX Annual Technical Conference (USENIX ATC 20).* USENIX Association, 617–631. https://www.usenix.org/conference/atc20/presentation/raghavan

[19] Radu Rugina and Martin C. Rinard. 1999. Automatic Parallelization of Divide and Conquer Algorithms. In *Proceedings of the 1999 ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPOPP'99), Atlanta, Georgia, USA, May 4-6, 1999,* Marc Snir and Andrew A. Chien (Eds.). ACM, 72–83.

[20] Wei Shen. 2019. A Cross-platform Command-line Tool for Executing Jobs in Parallel. https://github.com/shenwei356/rush.

[21] Calvin Smith and Aws Albarghouthi. 2016. MapReduce Program Synthesis. In *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Santa Barbara, CA, USA) *(PLDI '16).* Association for Computing Machinery, New York, NY, USA, 326–340. https://doi.org/10.1145/2908080.2908102

[22] Diomidis Spinellis and Marios Fragkoulis. 2017. Extending Unix Pipelines to DAGs. *IEEE Trans. Comput.* 66, 9 (2017), 1547–1561.

[23] Ole Tange. 2011. GNU Parallel—The Command-Line Power Tool. *;login: The USENIX Magazine* 36, 1 (Feb 2011), 42–47. https://doi.org/10.5281/zenodo.16303

[24] Dave Taylor. 2004. *Wicked Cool Shell Scripts: 101 Scripts for Linux, Mac OS X, and Unix Systems.* No Starch Press.

[25] Eleftheria Tsaliki and Diomidis Spinellis. 2021. The real statistics of buses in Athens. https://insidestory.gr/article/noymera-leoforeia-athinas?token=0MFVISB8N6.

[26] Nikos Vasilakis, Konstantinos Kallas, Konstantinos Mamouras, Achilles Benetopoulos, and Lazar Cvetković. 2021. PaSh: Light-Touch Data-Parallel Shell Processing. In *Proceedings of the Sixteenth European Conference on Computer Systems.* Association for Computing Machinery, New York, NY, USA, 49–66. https://doi.org/10.1145/3447786.3456228

[27] Andy B Yoo, Morris A Jette, and Mark Grondona. 2003. Slurm: Simple linux utility for resource management. In *Workshop on Job Scheduling Strategies for Parallel Processing.* Springer, 44–60.

$$\frac{i_1 = \text{strToInt } y_1 \qquad i_2 = \text{strToInt } y_2}{\textbf{add } y_1 \ y_2 \Longrightarrow_e \text{intToStr } (i_1 + i_2)}$$

$$\frac{}{\textbf{concat } y_1 \ y_2 \Longrightarrow_e y_1 \text{ ++ } y_2}$$

$$\frac{}{\textbf{first } y_1 \ y_2 \Longrightarrow_e y_1}$$

$$\frac{}{\textbf{second } y_1 \ y_2 \Longrightarrow_e y_2}$$

$$\frac{b \ (\text{delFront } d \ y_1) \ (\text{delFront } d \ y_2) \Longrightarrow_e v}{(\textbf{front } d \ b) \ y_1 \ y_2 \Longrightarrow_e d \text{ ++ } v}$$

$$\frac{b \ (\text{delBack } d \ y_1) \ (\text{delBack } d \ y_2) \Longrightarrow_e v}{(\textbf{back } d \ b) \ y_1 \ y_2 \Longrightarrow_e v \text{ ++ } d}$$

$$\frac{\begin{array}{c} h_1, t_1 = \text{splitFirst } d \ y_1 \qquad h_2, t_2 = \text{splitFirst } d \ y_2 \\ t_1 \neq \text{nil} \qquad t_2 \neq \text{nil} \qquad d \notin t_1 \qquad d \notin t_2 \\ b \ h_1 \ h_2 \Longrightarrow_e v \qquad b \ t_1 \ t_2 \Longrightarrow_e v' \end{array}}{(\textbf{fuse } d \ b) \ y_1 \ y_2 \Longrightarrow_e v \text{ ++ } d \text{ ++ } v'}$$

$$\frac{\begin{array}{c} h_1, t_1 = \text{splitFirst } d \ y_1 \qquad h_2, t_2 = \text{splitFirst } d \ y_2 \\ t_1 \neq \text{nil} \qquad t_2 \neq \text{nil} \qquad d \in t_1 \qquad d \in t_2 \\ b \ h_1 \ h_2 \Longrightarrow_e v \qquad (\textbf{fuse } d \ b) \ t_1 \ t_2 \Longrightarrow_e v' \end{array}}{(\textbf{fuse } d \ b) \ y_1 \ y_2 \Longrightarrow_e v \text{ ++ } d \text{ ++ } v'}$$

$$\frac{}{\textbf{rerun}_f \ y_1 \ y_2 \Longrightarrow_e f \ (y_1 \text{ ++ } y_2)}$$

$$\frac{v = (\text{unixMerge <flags>}) \ y_1 \ y_2}{(\textbf{merge <flags>}) \ y_1 \ y_2 \Longrightarrow_e v}$$

$$\frac{y_1 = \text{`}\backslash n\text{`} \text{ or } y_2 = \text{`}\backslash n\text{`}}{(\textbf{stitch } b) \ y_1 \ y_2 \Longrightarrow_e y_1 \text{ ++ } y_2}$$

$$\frac{y_1', l_1 = \text{splitLastLine } y_1 \qquad l_2, y_2' = \text{splitFirstLine } y_2 \qquad l_1 \neq l_2}{(\textbf{stitch } b) \ y_1 \ y_2 \Longrightarrow_e y_1 \text{ ++ } y_2}$$

$$\frac{y_1', l_1 = \text{splitLastLine } y_1 \quad l_2, y_2' = \text{splitFirstLine } y_2 \quad l_1 = l_2 \quad b \ l_1 \ l_2 \Longrightarrow_e v}{(\textbf{stitch } b) \ y_1 \ y_2 \Longrightarrow_e y_1' \text{ ++ `}\backslash n\text{` ++ } v \text{ ++ `}\backslash n\text{` ++ } y_2'}$$

$$\frac{\begin{array}{c} y_1', l_1 = \text{splitLastLine } v_1 \qquad h_1, t_1 = \text{splitFirst } d \ (\text{delPad } l_1) \\ l_2, y_2' = \text{splitFirstLine } v_2 \qquad h_2, t_2 = \text{splitFirst } d \ (\text{delPad } l_2) \qquad t_1 \neq t_2 \end{array}}{(\textbf{stitch2 } d \ b_1 \ b_2) \ v_1 \ v_2 \Longrightarrow_e y_1 \text{ ++ } y_2}$$

$$\frac{\begin{array}{c} y_1', l_1 = \text{splitLastLine } v_1 \qquad h_1, t_1 = \text{splitFirst } d \ (\text{delPad } l_1) \\ l_2, y_2' = \text{splitFirstLine } v_2 \qquad h_2, t_2 = \text{splitFirst } d \ (\text{delPad } l_2) \\ t_1 = t_2 \qquad b_1 \ h_1 \ h_2 \Longrightarrow_e h \qquad b_2 \ t_1 \ t_2 \Longrightarrow_e t \qquad v = \text{addPad } (h \text{ ++ } d \text{ ++ } t) \end{array}}{(\textbf{stitch2 } d \ b_1 \ b_2) \ v_1 \ v_2 \Longrightarrow_e y_1' \text{ ++ `}\backslash n\text{` ++ } v \text{ ++ `}\backslash n\text{` ++ } y_2'}$$

$$\frac{}{(\textbf{helper } d \ b) \ h_1 \ \text{nil} \Longrightarrow_e \text{nil}}$$

$$\frac{l_2, y_2' = \text{splitFirstLine } y_2 \qquad l_2 = \text{nil} \qquad (\textbf{helper } d \ b) \ h_1 \ y_2' \Longrightarrow_e v}{(\textbf{helper } d \ b) \ h_1 \ y_2 \Longrightarrow_e \text{`}\backslash n\text{` ++ } v}$$

$$\frac{\begin{array}{c} l_2, y_2' = \text{splitFirstLine } y_2 \qquad h_2, t_2 = \text{splitFirst } d \ (\text{delPad } l_2) \\ b \ h_1 \ h_2 \Longrightarrow_e h \qquad v = \text{addPad } (h \text{ ++ } d \text{ ++ } t_2) \qquad (\textbf{helper } d \ b) \ h_1 \ y_2' \Longrightarrow_e v' \end{array}}{(\textbf{helper } d \ b) \ h_1 \ y_2 \Longrightarrow_e v \text{ ++ `}\backslash n\text{` ++ } v'}$$

$$\frac{\begin{array}{c} y_1', l_1 = \text{splitLastNonemptyLine } y_1 \qquad h_1, t_1 = \text{splitFirst } d \ (\text{delPad } l_1) \\ (\textbf{helper } d \ b) \ h_1 \ y_2 \Longrightarrow_e v \end{array}}{(\textbf{offset } d \ b) \ y_1 \ y_2 \Longrightarrow_e y_1 \text{ ++ } v}$$

$$b, b_1, b_2 \in \text{RecOp} \qquad d \in \text{Delim} \qquad y_1, y_2, y_1', y_2', v, v', v_1, v_2, h, h_1, h_2, t, t_1, t_2, l_1, l_2 \in \text{String} \qquad i_1, i_2 \in \text{Int}$$

**Figure 6. DSL Semantics.** The semantics of synthesizable combiners $g \in \text{Combiner}_f$ for command $f$. A combiner accepts two strings $y_1, y_2$ that are the outputs from two executions of $f$. A plausible combiner $g$ for $f$ must satisfy $g \ y_1 \ y_2 \Longrightarrow_e f(x_1 \text{ ++ } x_2)$ for all of the observed input pairs $\langle x_1, x_2 \rangle$, where $y_1 = f(x_1)$ and $y_2 = f(x_2)$.

We present the combiner DSL semantics in Appendix A, a correctness result in Appendix B, performance results in Appendix C, and combiner synthesis results in Appendix D.

## A  Appendix: DSL Semantics

Figure 6 presents the big-step execution semantics for the DSL. The transition function $\Rightarrow$ maps a DSL expression to its output value. strToInt converts a string into an integer. intToStr converts an integer into a string. ++ concatenates two strings. unixMerge takes a comparator flag and two strings, then uses the flag to execute the "sort -m" command to merge the two strings. delFront and delBack each takes a delimiter and a string. delFront removes the specified delimiter from the beginning of the string, while delBack removes the delimiter at the end of the stream. nil denotes an empty string. splitFirst, splitLast, and splitLastNonempty each takes a delimiter and a string. splitFirst splits the string into elements separated by the delimiter, then returns the first element as the first output. It connects the remaining elements using the delimiter as the second output. splitLast likewise splits the string with the delimiter, then returns the last element as the second output and returns the remaining substring as the first output. splitLastNonempty splits the string with the delimiter, then returns the last nonempty element. delPad removes leading spaces from a string, then returns the number of removed spaces as the first output and returns the remaining substring as the second output. calcPad takes an integer and two strings, where the integer denotes the number of spaces that pad the first string. It returns the padding needed for the second string. addPad inserts padding before a string.

# B   Appendix: Conditions for Synthesizing Correct Combiners

**Definition B.1.** For $g \in \text{Combiner}_f$, $L(g)$ denotes the set of legal strings for which $g$ is defined.

$$L(\textbf{add}) = [\text{`0`} - \text{`9`}]^+$$

$$L(\textbf{concat}) = \text{String}$$

$$L(\textbf{first}) = \text{String}$$

$$L(\textbf{second}) = \text{String}$$

$$L(\textbf{front } d \text{ } b) = \{d \mathbin{++} y \mid y \in L(b)\}$$

$$L(\textbf{back } d \text{ } b) = \{y \mathbin{++} d \mid y \in L(b)\}$$

$$L(\textbf{fuse } d \text{ } b) = \{y_1 \mathbin{++} d \mathbin{++} y_2 \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} y_k$$
$$\mid y_1 \neq \text{nil}, y_k \neq \text{nil, and } y_i \in L(b) \text{ and } d \notin y_i$$
$$\text{for all } i = 1, \ldots, k, \text{ where } k \geq 2\}$$

$$L(\textbf{stitch } b) = \{y_1 \mathbin{++} \text{`} \backslash n \text{`} \mathbin{++} \ldots \mathbin{++} y_k \mathbin{++} \text{`} \backslash n \text{`}$$
$$\mid y_i \in L(b) \text{ and `} \backslash n \text{`} \notin y_i \text{ for all } i = 1, \ldots, k$$
$$\text{where } k \geq 1\}$$
$$\cup \{\text{`} \backslash n \text{`}\}$$

$$L(\textbf{stitch2 } d \text{ } b_1 \text{ } b_2) = \{y_1 \mathbin{++} \text{`} \backslash n \text{`} \mathbin{++} \ldots \mathbin{++} y_k \mathbin{++} \text{`} \backslash n \text{`}$$
$$\mid y_i = p \mathbin{++} h_i \mathbin{++} d \mathbin{++} t_i \text{ and `} \backslash n \text{`} \notin y_i$$
$$\text{for all } i = 1, \ldots, k, \text{ where } k \geq 1, p \in [\text{` `}^+ \mid \text{`} \backslash t \text{`}],$$
$$h_i \in L(b_1), d \notin h_i, \text{ and } t_i \in L(b_2)\}$$
$$\cup \{\text{`} \backslash n \text{`}\}$$

$$L(\textbf{offset } d \text{ } b) = \{y_1 \mathbin{++} \text{`} \backslash n \text{`} \mathbin{++} \ldots \mathbin{++} y_k \mathbin{++} \text{`} \backslash n \text{`}$$
$$\mid y_i \in \{\text{nil}, (p \mathbin{++} h_i \mathbin{++} d \mathbin{++} t_i)\} \text{ and `} \backslash n \text{`} \notin y_i$$
$$\text{for all } i = 1, \ldots, k, \text{ where } k \geq 1, p \in [\text{` `}^+ \mid \text{`} \backslash t \text{`}],$$
$$h_i \in L(b), d \notin h_i, \text{ and } t_i \in \text{String}\}$$

$$L(\textbf{rerun}_f) = \{\text{legal inputs for } f\}$$

$$L(\textbf{merge <flags>}) = \{\text{legal inputs for (unixMerge <flags>)}\}$$

For any $y_1, y_2 \in L(g)$, the evaluation $g$ $y_1$ $y_2 \Longrightarrow_e v$ succeeds for some $v \in \text{String}$.

**Definition B.2.** A stream is a string that ends with a newline character `$\backslash n$`, $\text{Stream} = \{x \mathbin{++} \text{`} \backslash n \text{`} \mid x \in \text{String}\}$.

**Definition B.3.** A command $f : \text{Stream} \to \text{Stream}$ is a function that takes a stream as input and produces a stream as output.[6]

**Definition B.4.** An input pair $\langle x_1, x_2 \rangle$ consists of two strings $x_1, x_2 \in \text{String}$. An output tuple $\langle y_1, y_2, y_{12} \rangle$ consists of three strings $y_1, y_2, y_{12} \in \text{String}$.

**Definition B.5.** An input stream pair $\langle x_1, x_2 \rangle$ consists of two streams $x_1, x_2 \in \text{Stream}$. An observation $\langle y_1, y_2, y_{12} \rangle$ consists of three streams $y_1, y_2, y_{12} \in \text{Stream}$.

**Definition B.6.** Executing command $f$ with an input stream pair $\langle x_1, x_2 \rangle$ produces the observation $\langle f(x_1), f(x_2), f(x_1 \mathbin{++} x_2) \rangle$. For a set of input stream pairs $X$, $f(X)$ denotes the set of observations obtained from executing $f$ with $X$, $f(X) = \{\langle f(x_1), f(x_2), f(x_1 \mathbin{++} x_2) \rangle \mid \langle x_1, x_2 \rangle \in X\}$.

**Definition B.7.** For $g_1, g_2 \in \text{Combiner}_f$, $g_1$ and $g_2$ are equivalent by intersection, denoted as $g_1 \equiv_\cap g_2$, if for all $y_1, y_2 \in L(g_1) \cap L(g_2)$, $g_1$ $y_1$ $y_2 \Longrightarrow_e v$ and $g_2$ $y_1$ $y_2 \Longrightarrow_e v$ for some $v$.

*Example* 1. For all $d \in \text{Delim}$, we have $(\textbf{front } d \textbf{ concat}) \equiv_\cap (\textbf{back } d \textbf{ concat})$ and $(\textbf{stitch2 } d \textbf{ first first}) \equiv_\cap (\textbf{stitch first})$.

**Definition B.8.** A combiner $g \in \text{Combiner}_f$ is plausible for output tuples $Y$, denoted as $P(g, Y)$, if $y_1, y_2 \in L(g)$ and $g$ $y_1$ $y_2 \Longrightarrow_e y_{12}$ for all $\langle y_1, y_2, y_{12} \rangle \in Y$.

---

[6]Although this paper focuses on commands whose outputs terminate with newlines, the KᴜᴍQᴜᴀᴛ algorithm applies also to commands whose outputs do not terminate with newlines.

**Table 2.** Representative combiners

| Combiner $g$ | Stage | Conditions for $E(g, Y)$ to be true |
|---|---|---|
| $g_a = \textbf{add}$ | RecOp | The following conditions hold: (1) There exists $\langle y_1, y_2, y_{12}\rangle \in Y$ such that $y_1 \notin$ '0'$^+$. (2) There exists $\langle y_1, y_2, y_{12}\rangle \in Y$ such that $y_2 \notin$ '0'$^+$. |
| $g_c = \textbf{concat}$ | RecOp | The following conditions hold: (1) There exists $\langle y_1, y_2, y_{12}\rangle \in Y$ such that $y_1 \neq$ nil. (2) There exists $\langle y_1, y_2, y_{12}\rangle \in Y$ such that $y_2 \neq$ nil. |
| $g_f = \textbf{first}$ | RecOp | The following conditions hold: (1) There exists $\langle y_1, y_2, y_{12}\rangle \in Y$ such that $y_1 \neq y_2$. (2) There exists $\langle y_1, y_2, y_{12}\rangle \in Y$ and $c \in y_2$ such that $c \notin$ Delim $\cup$ {'0'}. |
| $g_s = \textbf{second}$ | RecOp | The following conditions hold: (1) There exists $\langle y_1, y_2, y_{12}\rangle \in Y$ such that $y_1 \neq y_2$. (2) There exists $\langle y_1, y_2, y_{12}\rangle \in Y$ and $c \in y_1$ such that $c \notin$ Delim $\cup$ {'0'}. |
| $g_{ba} = (\textbf{back } d \textbf{ add})$ | RecOp | $E(g_a, Y')$ where $Y' = \{\langle y_1, y_2, y_{12}\rangle \mid \langle (y_1 {+}{+} d), (y_2 {+}{+} d), (y_{12} {+}{+} d)\rangle \in Y\}$. |
| $g_{fa} = (\textbf{fuse } d \textbf{ add})$ | RecOp | $E(g_a, Y')$ where $Y' = \{\langle y_{1,i}, y_{2,i}, y_{12,i}\rangle \mid \langle (y_{1,1} {+}{+} d {+}{+} \dots {+}{+} d {+}{+} y_{1,n}), (y_{2,1} {+}{+} d {+}{+} \dots {+}{+} d {+}{+} y_{2,n}), (y_{12,1} {+}{+} d {+}{+} \dots {+}{+} d {+}{+} y_{12,n})\rangle \in Y$, where $d \notin y_j^1$ and $d \notin y_j^2$ for all $j = 1, \dots, n$, and $i \in \{1, \dots, n\}\}$. |
| $g_{bfa} = (\textbf{back } d_1 (\textbf{fuse } d_2 \textbf{ add}))$ | RecOp | $E(g_{fa}, Y')$ where $Y' = \{\langle y_1, y_2, y_{12}\rangle \mid \langle (y_1 {+}{+} d), (y_2 {+}{+} d), (y_{12} {+}{+} d)\rangle \in Y\}$. |
| $g_{fbfa} = (\textbf{front } d_1 (\textbf{back } d_2 (\textbf{fuse } d_3 \textbf{ add})))$ | RecOp | $E(g_{bfa}, Y')$ where $Y' = \{\langle y_1, y_2, y_{12}\rangle \mid \langle (d {+}{+} y_1), (d {+}{+} y_2), (d {+}{+} y_{12})\rangle \in Y\}$. |
| $g_{fc} = (\textbf{front } d \textbf{ concat})$ | RecOp | $E(g_c, Y')$ where $Y' = \{\langle y_1, y_2, y_{12}\rangle \mid \langle (d {+}{+} y_1), (d {+}{+} y_2), (d {+}{+} y_{12})\rangle \in Y\}$. |
| $g_{sf} = (\textbf{stitch first})$ | StructOp | The following conditions hold: (1) There exists $\langle y_1, y_2, y_{12}\rangle \in Y$ such that (splitLastLine $y_1$) $= (y_1', l)$ and (splitFirstLine $y_2$) $= (l, y_2')$ where (firstChar (delPad $l$)) $\notin$ Delim $\cup$ {'0'} and (lastChar $l$) $\notin$ Delim $\cup$ {'0'}. (2) If $Y \subseteq L(\textbf{stitch2 } d \; b_1 \; b_2)$ for some $d \in$ Delim and $b_1, b_2 \in$ RecOp, then there exists $\langle y_1, y_2, y_{12}\rangle \in Y$ such that $h_1 \neq h_2$, where (splitLastLine $y_1$) $= (y_1', l_1)$, (splitFirstLine $y_2$) $= (l_2, y_2')$, (splitFirst $d$ (delPad $l_1$)) $= (h_1, t)$, and (splitFirst $d$ (delPad $l_2$)) $= (h_2, t)$. |
| $g_{saf} = (\textbf{stitch2 } d \textbf{ add first})$ | StructOp | There exists $\langle y_1, y_2, y_{12}\rangle \in Y$ such that (splitLastLine $y_1$) $= (y_1', l)$ and (splitFirstLine $y_2$) $= (l, y_2')$ where (firstChar (delPad $l$)) $\notin$ Delim $\cup$ {'0'} and (lastChar $l$) $\notin$ Delim $\cup$ {'0'}. |
| $g_{oa} = (\textbf{offset } d \textbf{ add})$ | StructOp | The following conditions hold: (1) There exists $\langle y_1, y_2, y_{12}\rangle \in Y$ such that (splitLastLine $y_1$) $= (y_1', l_1)$, (splitFirstLine $y_2$) $= (l_2, y_2')$, and (splitFirstLine $y_2'$) $= (l_2', y_2'')$, where (firstChar (delPad $l_1$)) $\notin$ Delim $\cup$ {'0'}, $l_2 \neq$ nil, and $l_2' \neq$ nil. (2) $E(g_a, Y')$ where $Y' = \{\langle h_1, h_2, y_{12}'\rangle \mid \langle y_1, y_2, y_{12}\rangle \in Y,$ (splitLastLine $y_1$) $= (y_1', l_1)$, (splitFirstLine $y_2$) $= (l_2, y_2')$, (splitFirst $d$ (delPad $l_1$)) $= (h_1, t_1)$, and (splitFirst $d$ (delPad $l_2$)) $= (h_2, t_2)\}$. |

**Definition B.9.** A combiner $g \in$ Combiner$_f$ is correct for command $f$ if $P(g, f(X))$ holds for all input stream pairs $X$.

*Remark.* Note that this definition does not require $P(g, Y)$ for all sets of output tuples $Y$.

**Definition B.10.** $C(d, y)$ denotes the number of times that $d \in$ Delim occurs in $y \in$ String. We write $d \in y$ when $C(d, y) > 0$ and write $d \notin y$ when $C(d, y) = 0$.

**Definition B.11.** We define two sets of representative combiners for command $f$, $G_{rec} = \{g_a, g_c, g_f, g_s, g_{ba}, g_{fa}, g_{bfa}, g_{fbfa}, g_{fc}\} \subset$ RecOp and $G_{struct} = \{g_{sf}, g_{saf}, g_{oa}\} \subset$ StructOp, whose elements are defined as follows:

$$g_a = \textbf{add} \in \text{RecOp}$$
$$g_c = \textbf{concat} \in \text{RecOp}$$
$$g_f = \textbf{first} \in \text{RecOp}$$
$$g_s = \textbf{second} \in \text{RecOp}$$
$$g_{ba} = (\textbf{back } d \textbf{ add}) \in \text{RecOp}$$
$$g_{fa} = (\textbf{fuse } d \textbf{ add}) \in \text{RecOp}$$
$$g_{bfa} = (\textbf{back } d_1 (\textbf{fuse } d_2 \textbf{ add})) \in \text{RecOp}$$
$$g_{fbfa} = (\textbf{front } d_1 (\textbf{back } d_2 (\textbf{fuse } d_3 \textbf{ add}))) \in \text{RecOp}$$
$$g_{fc} = (\textbf{front } d \textbf{ concat}) \in \text{RecOp}$$
$$g_{sf} = (\textbf{stitch first}) \in \text{StructOp}$$
$$g_{saf} = (\textbf{stitch2 } d \textbf{ add first}) \in \text{StructOp}$$
$$g_{oa} = (\textbf{offset } d \textbf{ add}) \in \text{StructOp}$$

**Definition B.12.** For combiner $g \in G_{rec} \cup G_{struct}$ and any set of output tuples $Y$, $E(g, Y)$ denotes a conservative predicate that is true only if $Y$ is sufficient for eliminating incorrect candidates when the correct combiner is $g$. We define these predicates in Table 2.

**Definition B.13.** For any set of output tuples $Y$, $E_{\text{rec}}(Y)$ denotes a conservative predicate that is true only if $Y$ is sufficient for eliminating incorrect candidates when the correct combiner $g \in G_{\text{rec}}$. $E_{\text{rec}}(Y)$ is true if and only if the following conditions hold:

- There exists $\langle y_1, y_2, y_{12} \rangle \in Y$ such that $y_1 \neq y_2$.
- There exists $\langle y_1, y_2, y_{12} \rangle \in Y$ and $c \in y_1$ such that $c \notin \text{Delim} \cup \{\text{'0'}\}$.
- There exists $\langle y_1, y_2, y_{12} \rangle \in Y$ and $c \in y_2$ such that $c \notin \text{Delim} \cup \{\text{'0'}\}$.

**Definition B.14.** For any set of output tuples $Y$, $T(Y)$ denotes a predicate that is true only if $Y$ is interpretable as a table. $T(Y)$ is true if and only if there exists $p \in [\text{' '+} \mid \text{'\textbackslash t'}]$ and $d \in \text{Delim}$ such that for all $\langle y_1, y_2, y_{12} \rangle \in Y$, each line in $y_1, y_2, y_{12}$ is either nil or of the form $(p \mathbin{++} h \mathbin{++} d \mathbin{++} t)$ for some $h, t \in \text{String}$.

**Definition B.15.** For any set of output tuples $Y$, $E_{\text{struct}}(Y)$ denotes a conservative predicate that is true only if $Y$ is sufficient for eliminating incorrect candidates when the correct combiner $g \in G_{\text{struct}}$. $E_{\text{struct}}(Y)$ is true if and only if the following conditions hold:

- There exists $\langle y_1, y_2, y_{12} \rangle \in Y$ such that $(\text{splitLastLine } y_1) = (y_1', l)$, $(\text{splitFirstLine } y_2) = (l, y_2')$, and $(\text{splitFirstLine } y_2') = (l_2', y_2'')$, where $(\text{firstChar } (\text{delPad } l)) \notin \text{Delim} \cup \{\text{'0'}\}$, $(\text{lastChar } l) \notin \text{Delim} \cup \{\text{'0'}\}$, and $l_2' \neq \text{nil}$.
- If $T(Y)$ then $E_{\text{rec}}(Y')$, where $Y' = \{\langle h_1, h_2, y_{12}' \rangle \mid \langle y_1, y_2, y_{12} \rangle \in Y, (\text{splitLastLine } y_1) = (y_1', l_1), (\text{splitFirstLine } y_2) = (l_2, y_2'), (\text{splitFirst } d \; (\text{delPad } l_1)) = (h_1, t), \text{ and } (\text{splitFirst } d \; (\text{delPad } l_2)) = (h_2, t)\}$.

**Definition B.16.** The size of a combiner $g \in \text{Combiner}_f$ is denoted as $|g|$ and defined as two (each combiner operates on two arguments) plus the number of times that the AST of $g$ applies a production to expand a "RecOp", "StructOp", or "RunOp$_f$" symbol.

*Example* 2. We have $|g_a| = 3$, $|g_{\text{fbfa}}| = 6$, and $|g_{\text{saf}}| = 5$.

**Definition B.17.** For a command $f$, integer $k$, and set of output tuples $Y$, the set of plausible combiners $P_k(Y) = \{g \in \text{Combiner}_f \mid |g| \leq k \text{ and } P(g, Y)\}$.

**Proposition B.1.** If combiner $g \in \text{Combiner}_f$ is correct for command $f$, then $y_1, y_2 \in L(g)$ holds for all $\langle y_1, y_2, y_{12} \rangle \in f(X)$ where $X$ is any set of input stream pairs.

**Proposition B.2.** For any combiner $g \in G_{\text{rec}}$ and set of output tuples $Y$, if $P(g, Y)$ and $E_{\text{rec}}(Y)$ then $E(g, Y)$.

*Proof.* By induction on the derivation of $g$. □

**Proposition B.3.** For any $d \in \text{Delim}$ and $b_1, b_2 \in \text{RecOp}$, $Y \subseteq L(\textbf{stitch2 } d \; b_1 \; b_2)$ implies $T(Y)$. For any $d \in \text{Delim}$ and $b \in \text{RecOp}$, $Y \subseteq L(\textbf{offset } d \; b)$ implies $T(Y)$.

**Proposition B.4.** For any combiner $g \in G_{\text{struct}}$ and set of output tuples $Y$, if $P(g, Y)$ and $E_{\text{struct}}(Y)$ then $E(g, Y)$.

**Proposition B.5.** For any integers $k_1, k_2$ such that $0 < k_1 < k_2$, we have $P_{k_1}(Y) \subseteq P_{k_2}(Y)$ for all sets of output tuples $Y$.

**Proposition B.6.** If combiner $g \in \text{Combiner}_f$ is correct for command $f$, then $g \in P_{|g|}(f(X))$ for all sets of input stream pairs $X$.

**Proposition B.7.** For any integer $k \geq 6$ and set of input stream pairs $X$, we have:

- If combiner $g \in G_{\text{rec}}$ is correct for command $f$, then $g \in P_k(f(X)) \cap \text{RecOp}$.
- If combiner $g \in G_{\text{rec}}$ is correct for command $f$, then $g \in P_k(f(X)) \cap \text{StructOp}$.

**Lemma B.1.** For any $g \in \text{RecOp}$, $y_1, y_2 \in \text{String}$, and $d \in \text{Delim}$, if $g \; y_1 \; y_2 \Longrightarrow_e v$ and $d \notin y_1$ and $d \notin y_2$ then $d \notin v$.

*Proof.* By induction on the derivation of $g$. □

**Lemma B.2.** For any $g \in \text{RecOp}$, $y_1, y_2 \in \text{String}$, and $z \in \text{String}$, if $g \; y_1 \; y_2 \Longrightarrow_e v$ and $z \neq \text{nil}$ then $v \neq y_1 \mathbin{++} z \mathbin{++} y_2$.

*Proof.* By induction on the derivation of $g$, where the case of "$g = \textbf{fuse } d \; b$" uses Lemma B.1. □

**Lemma B.3.** For any $d \in \text{Delim}$, $b \in \text{RecOp}$, $g = \textbf{fuse } d \; b$, $y_1, y_2 \in L(g)$, and $y_{12} \in \text{String}$ such that $g \; y_1 \; y_2 \Longrightarrow_e y_{12}$, we have $C(d, y_1) = C(d, y_2) = C(d, y_{12})$.

*Proof.* By Definition B.8 and Definition B.1, $y_1, y_2 \in L(g') = \{y_1 \mathbin{++} d \mathbin{++} y_2 \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} y_k \mid y_i \in L(b) \text{ and } d \notin y_i \text{ for all } i = 1, \ldots, k, \text{ where } k \geq 2\}$. Let $k = C(d, y_1) + 1$. By Figure 6, $C(d, y_2) = C(d, y_1) = k-1$. Let $y_1 = y_1^1 \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} y_k^1$ and $y_2 = y_1^2 \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} y_k^2$ where $y_i^1, y_i^2 \in L(b)$, $d \notin y_i^1$, and $d \notin y_i^2$ for all $i = 1, \ldots, k$. By ??, $d \in \text{Delim}$ and $b \in \text{RecOp}$. By Figure 6, there exists $v_1, \ldots, v_k \in \text{String}$ such that $y_{12} = v_1 \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} v_k$ and $b \; y_i^1 \; y_i^2 \Longrightarrow_e v_i$ for all $i = 1, \ldots, k$. By Lemma B.1, $d \notin v_i$ for all $i = 1, \ldots, k$. Hence $C(d, y_{12}) = k - 1$. □

**Lemma B.4.** For any $d \in$ Delim, $g \in$ RecOp, and $y_1, y_2, y_{12} \in$ String such that $g\ y_1\ y_2 \Longrightarrow_e y_{12}$, we have $C(d, y_{12}) \leq C(d, y_1) + C(d, y_2)$.

*Proof.* By induction on the derivation of $g$. □

**Theorem 6.** For any combiner $g \in G_{\text{rec}}$, set of output tuples $Y$ such that $P(g, Y)$ and $E(g, Y)$, and $g' \in$ RecOp, we have $P(g', Y)$ implies $g' \equiv_\cap g$.

*Proof.* The proof is by induction on the derivation of $g$.

<u>Case 1:</u> $g = g_a$. The proof performs case analysis of the values of $g'$. By Table 2, there exists $\langle y_1, y_2, y_{12} \rangle \in Y$. By Definition B.8 and Definition B.1, $y_1, y_2 \in L(g) = ['0' - '9']^+$. By Definition B.8, $g\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = \text{intToStr}((\text{strToInt}\ y_1) + (\text{strToInt}\ y_2))$.

> <u>Case 1.1:</u> $g' = \textbf{add} = g$. By Definition B.7, $g' \equiv_\cap g$.
>
> <u>Case 1.2:</u> $g' = \textbf{concat}$. We show that $P(g', Y)$ never holds in this case. If $y_1 \notin '0'^+$, we have $(\text{strToInt}\ y_1) > 0$. Hence $\text{strToInt}(y_1 ++ y_2) \geq 10 \cdot (\text{strToInt}\ y_1) + (\text{strToInt}\ y_2) > (\text{strToInt}\ y_1) + (\text{strToInt}\ y_2)$. If $y_1 \in '0'^+$, the string $(y_1 ++ y_2)$ contains leading '0' characters that are absent in the results of intToStr. Either case, we have $\text{intToStr}((\text{strToInt}\ y_1) + (\text{strToInt}\ y_2)) \neq y_1 ++ y_2$. Assume the opposite that $P(g', Y)$. By Definition B.8, $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = y_1 ++ y_2$. We have the desired contradiction.
>
> <u>Case 1.3:</u> $g' = \textbf{first}$. We show that $P(g', Y)$ never holds in this case. By Table 2, there exists $\langle y_1, y_2, y_{12} \rangle \in Y$ such that $y_2 \notin '0'^+$. Hence $(\text{strToInt}\ y_2) > 0$. We have $\text{intToStr}((\text{strToInt}\ y_1) + (\text{strToInt}\ y_2)) \neq y_1$. Assume the opposite that $P(g', Y)$. By Definition B.8, $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = y_1$. We have the desired contradiction.
>
> <u>Case 1.4:</u> $g' = \textbf{second}$. The proof is similar to the proof of Case 1.3.
>
> <u>Case 1.5:</u> $g' = \textbf{front}\ d\ b$. We show that $P(g', Y)$ never holds in this case. Assume the opposite that $P(g', Y)$. By Definition B.8 and Definition B.1, $y_1, y_2 \in L(g') = \{d ++ y \mid y \in L(b)\}$. By Figure 3, $d \in$ Delim and $b \in$ RecOp. Thus, $d \notin ['0' - '9']$ and $L(g) \cap L(g') = \emptyset$. We have the desired contradiction.
>
> <u>Case 1.6:</u> $g' = \textbf{back}\ d\ b$. The proof is similar to the proof of Case 1.5.
>
> <u>Case 1.7:</u> $g' = \textbf{fuse}\ d\ b$. The proof is similar to the proof of Case 1.5.

<u>Case 2:</u> $g = g_c$. The proof performs case analysis of the values of $g'$. By Table 2, there exists $\langle y_1, y_2, y_{12} \rangle \in Y$. By Definition B.8, $g\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = y_1 ++ y_2$.

> <u>Case 2.1:</u> $g' = \textbf{add}$. The proof is similar to the proof of Case 1.2.
>
> <u>Case 2.2:</u> $g' = \textbf{concat} = g$. By Definition B.7, $g' \equiv_\cap g$.
>
> <u>Case 2.3:</u> $g' = \textbf{first}$. We show that $P(g', Y)$ never holds in this case. By Table 2, there exists $\langle y_1, y_2, y_{12} \rangle \in Y$ such that $y_2 \neq$ nil. Hence $y_1 ++ y_2 \neq y_1$. Assume the opposite that $P(g', Y)$. By Definition B.8, $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = y_1$. We have the desired contradiction.
>
> <u>Case 2.4:</u> $g' = \textbf{second}$. The proof is similar to the proof of Case 2.3.
>
> <u>Case 2.5:</u> $g' = \textbf{front}\ d\ b$. We show that $P(g', Y)$ never holds in this case. Assume the opposite that $P(g', Y)$. By Definition B.8 and Definition B.1, $y_1, y_2 \in L(g') = \{d ++ y \mid y \in L(b)\}$. Let $y_1 = d ++ y_1', y_2 = d ++ y_2'$ where $y_1', y_2' \in L(b)$. By Figure 3, $d \in$ Delim and $b \in$ RecOp. By Definition B.8, $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, there exists $v \in$ String such that $y_{12} = d ++ v$ and $b\ y_1'\ y_2' \Longrightarrow_e v$. Since $y_{12} = y_1 ++ y_2 = d ++ y_1' ++ d ++ y_2'$, we have $v = y_1' ++ d ++ y_2'$. Since $d \neq$ nil, by Lemma B.2, $v \neq y_1' ++ d ++ y_2'$. We have the desired contradiction.
>
> <u>Case 2.6:</u> $g' = \textbf{back}\ d\ b$. The proof is similar to the proof of Case 2.5.
>
> <u>Case 2.7:</u> $g' = \textbf{fuse}\ d\ b$. We show that $P(g', Y)$ never holds in this case. Assume the opposite that $P(g', Y)$. By Definition B.8 and Definition B.1, $y_1, y_2 \in L(g') = \{y_1' ++ d ++ y_2' ++ d ++ \ldots ++ d ++ y_k' \mid y_1' \neq \text{nil}, y_k' \neq \text{nil}, \text{and } y_i' \in L(b) \text{ and } d \notin y_i' \text{ for all } i = 1, \ldots, k, \text{ where } k \geq 2\}$. Let $k = C(d, y_1) + 1$. By Figure 6, $C(d, y_1) \geq 1$ and $k \geq 2$. By Lemma B.3, $C(d, y_2) = C(d, y_{12}) = C(d, y_1) = k - 1$. Since $y_{12} = y_1 ++ y_2, C(d, y_{12}) = 2 \cdot (k - 1)$. Since $k \geq 2$, we have $k - 1 < 2 \cdot (k - 1)$. We have the desired contradiction.

<u>Case 3:</u> $g = g_f$. The proof is by induction on the derivation of $g'$. By Table 2, there exists $\langle y_1, y_2, y_{12} \rangle \in Y$. By Definition B.8, $g\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = y_1$.

> <u>Case 3.1:</u> $g' = \textbf{add}$. The proof is similar to the proof of Case 1.3.

<u>Case 3.2:</u> $g' = $ **concat**. The proof is similar to the proof of Case 2.3.

<u>Case 3.3:</u> $g' = $ **first** $= g$. By Definition B.7, $g' \equiv_\cap g$.

<u>Case 3.4:</u> $g' = $ **second**. We show that $P(g', Y)$ never holds in this case. By Table 2, there exists $\langle y_1, y_2, y_{12} \rangle \in Y$ such that $y_1 \neq y_2$. Assume the opposite that $P(g', Y)$. By Definition B.8, $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = y_2$. We have the desired contradiction.

<u>Case 3.5:</u> $g' = $ **front** $d\ b$. By Definition B.8 and Definition B.1, $y_1, y_2 \in L(g') = \{d ++ y \mid y \in L(b)\}$. Let $y_1 = d ++ y'_1$, $y_2 = d ++ y'_2$ where $y'_1, y'_2 \in L(b)$. By Figure 3, $d \in$ Delim and $b \in$ RecOp. By Definition B.8, $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, there exists $v \in$ String such that $y_{12} = d ++ v$ and $b\ y'_1\ y'_2 \Longrightarrow_e v$. Since $y_{12} = y_1 = d ++ y'_1$, we have $b\ y'_1\ y'_2 \Longrightarrow_e y'_1$.

Let $Y' = \{\langle y'_1, y'_2, y'_{12} \rangle \mid \langle (d ++ y'_1), (d ++ y'_2), (d ++ y'_{12}) \rangle \in Y\}$. We have $b\ y'_1\ y'_2 \Longrightarrow_e y'_{12}$ and $y'_{12} = y'_1$ for all $\langle y'_1, y'_2, y'_{12} \rangle \in Y'$. By Definition B.8, $P(b, Y')$ and $P(g, Y')$.

By Table 2, there exists $\langle y_1, y_2, y_{12} \rangle \in Y$ such that $y_1 \neq y_2$. Let $y_1 = d ++ y'_1$, $y_2 = d ++ y'_2$, $y_{12} = d ++ y'_{12}$. We have $y'_1 \neq y'_2$ and $\langle y'_1, y'_2, y'_{12} \rangle \in Y'$. By Table 2, there exists $\langle y_1, y_2, y_{12} \rangle \in Y$ and $c \in y_2$ such that $c \notin$ Delim $\cup \{$'0'$\}$. Let $y_2 = d ++ y'_2$. Since $d \in$ Delim, we have $c \in y'_2$.

By Table 2, $E(g, Y')$. By the induction hypothesis, $b \equiv_\cap g$.

For all $y'_1, y'_2 \in L(b) \cap L(g)$, by Definition B.7, $b\ y'_1\ y'_2 \Longrightarrow_e v$ and $g\ y'_1\ y'_2 \Longrightarrow_e v$ for some $v$. By Figure 6, $v = y'_1$. By Definition B.1, $L(g) =$ String and $L(b) \cap L(g) = L(b)$. Hence $b\ y'_1\ y'_2 \Longrightarrow_e y'_1$ for all $y'_1, y'_2 \in L(b)$.

For all $y_1, y_2 \in L(g')$, by Definition B.1, $y_1 = (d ++ y'_1)$ and $y_2 = (d ++ y'_2)$ for some $y'_1, y'_2 \in L(b)$. By Figure 6, $g'\ (d ++ y'_1)\ (d ++ y'_2) \Longrightarrow_e d ++ y'_1$. Hence $g'\ y_1\ y_2 \Longrightarrow_e y_1$ for all $y_1, y_2 \in L(g')$. Since $L(g) =$ String, $L(g') \cap L(g) = L(g')$. By Definition B.7, $g' \equiv_\cap g$.

<u>Case 3.6:</u> $g' = $ **back** $d\ b$. The proof is similar to the proof of Case 3.5.

<u>Case 3.7:</u> $g' = $ **fuse** $d\ b$. By Definition B.8 and Definition B.1, $y_1, y_2 \in L(g') = \{y'_1 ++ d ++ y'_2 ++ d ++ \dots ++ d ++ y'_k \mid y'_1 \neq$ nil, $y'_k \neq$ nil, and $y'_i \in L(b)$ and $d \notin y'_i$ for all $i = 1, \dots, k$, where $k \geq 2\}$. Let $y_1 = y_1^1 ++ d ++ \dots ++ d ++ y_k^1$ for some $y_i^1 \in L(b)$, where $d \notin y_i^1$ for all $i = 1, \dots, k$. By Lemma B.3, $y_2 = y_1^2 ++ d ++ \dots ++ d ++ y_k^2$ and $y_{12} = y_1^{12} ++ d ++ \dots ++ d ++ y_k^{12}$ for some $y_i^2 \in L(b)$, $y_i^{12} \in$ String, where $d \notin y_i^2$ and $d \notin y_i^{12}$ for all $i = 1, \dots, k$. By Definition B.8, $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $b\ y_i^1\ y_i^2 \Longrightarrow_e y_i^{12}$ for all $i = 1, \dots, k$. Since $y_{12} = y_1$, we have $y_i^{12} = y_i^1$ for all $i = 1, \dots, k$.

Let $Y' = \{\langle y_i^1, y_i^2, y_i^{12} \rangle \mid \langle (y_1^1 ++ d ++ \dots ++ d ++ y_k^1), (y_1^2 ++ d ++ \dots ++ d ++ y_k^2), (y_1^{12} ++ d ++ \dots ++ d ++ y_k^{12}) \rangle \in Y$, where $d \notin y_j^1$ and $d \notin y_j^2$ for all $j = 1, \dots, k$, and $i =\in \{1, \dots, k\}\}$. We have $b\ y_i^1\ y_i^2 \Longrightarrow_e y_i^{12}$ and $y_i^{12} = y_i^1$ for all $\langle y_i^1, y_i^2, y_i^{12} \rangle \in Y'$. By Definition B.8, $P(b, Y')$ and $P(g, Y')$.

The rest of the proof is similar to the proof of Case 3.5.

<u>Case 4:</u> $g = g_s$. The proof is similar to the proof of Case 3.

<u>Case 5:</u> $g = g_{ba}$. The proof is by induction on the derivation of $g'$. By Table 2, there exists $\langle y_1, y_2, y_{12} \rangle \in Y$. By Definition B.8 and Definition B.1, $y_1, y_2 \in L(g) = [$'0' $-$ '9'$]^+ d$.

<u>Case 5.1:</u> $g' = $ **add**. The proof is similar to the proof of Case 1.5.

<u>Case 5.2:</u> $g' = $ **concat**. We show that $P(g', Y)$ never holds in this case. By Definition B.8, $g\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} \in L(g)$. Assume the opposite that $P(g', Y)$. By Definition B.8, $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = y_1 ++ y_2$. Hence $y_{12} \in [$'0' $-$ '9'$]^+ d\ [$'0' $-$ '9'$]^+ d$. By Figure 3, $d \in$ Delim. Thus, $d \notin [$'0' $-$ '9'$]$ and $L(g) \cap [$'0' $-$ '9'$]^+ d\ [$'0' $-$ '9'$]^+ d = \emptyset$. We have the desired contradiction.

<u>Case 5.3:</u> $g' = $ **first**. We show that $P(g', Y)$ never holds in this case. Let $y_1 = y'_1 ++ d$, $y_2 = y'_2 ++ d$ where $y'_1, y'_2 \in [$'0' $-$ '9'$]^+$. By Definition B.8, $g\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $g_a\ y'_1\ y'_2 \Longrightarrow_e y'_1$. Assume the opposite that $P(g', Y)$. By Definition B.8, $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = y_1 = y'_1 ++ d$.

Let $Y' = \{\langle y'_1, y'_2, y'_{12} \rangle \mid \langle (y'_1 ++ d), (y'_2 ++ d), (y'_{12} ++ d) \rangle \in Y\}$. We have $g_a\ y'_1\ y'_2 \Longrightarrow_e y'_{12}$ and $y'_{12} = y'_1$ for all $\langle y'_1, y'_2, y'_{12} \rangle \in Y'$. By Definition B.8, $P(g_a, Y')$ and $P(g', Y')$.

By Table 2, $E(g_a, Y')$. We show the desired contradiction similar to the proof of Case 1.3.

<u>Case 5.4:</u> $g' = $ **second**. The proof is similar to the proof of Case 5.3.

<u>Case 5.5:</u> $g' = $ **front** $d'\ b$. The proof is similar to the proof of Case 1.5.

<u>Case 5.6:</u> $g' = $ **back** $d'\ b$ where $d' \neq d$. The proof is similar to the proof of Case 1.5.

<u>Case 5.7:</u> $g' = \textbf{back } d\ b$. Let $y_1 = y_1' \mathbin{++} d$, $y_2 = y_2' \mathbin{++} d$ where $y_1', y_2' \in [\text{'0'} - \text{'9'}]^+$. By Definition B.8, $g\ y_1\ y_2 \Longrightarrow_e y_{12}$ and $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, there exists $v \in \text{String}$ such that $y_{12} = v \mathbin{++} d$, $g_a\ y_1'\ y_2' \Longrightarrow_e v$, and $b\ y_1'\ y_2' \Longrightarrow_e v$.

Let $Y' = \{\langle y_1', y_2', y_{12}' \rangle \mid \langle (y_1' \mathbin{++} d), (y_2' \mathbin{++} d), (y_{12}' \mathbin{++} d) \rangle \in Y\}$. We have $g_a\ y_1'\ y_2' \Longrightarrow_e y_{12}'$ and $b\ y_1'\ y_2' \Longrightarrow_e y_{12}'$ for all $\langle y_1', y_2', y_{12}' \rangle \in Y'$. By Definition B.8, $P(g_a, Y')$ and $P(b, Y')$.

By Table 2, $E(g_a, Y')$. By the induction hypothesis, $b \equiv_\cap g_a$.

For all $y_1', y_2' \in L(b) \cap L(g_a)$, by Definition B.7, $b\ y_1'\ y_2' \Longrightarrow_e v$ and $g_a\ y_1'\ y_2' \Longrightarrow_e v$ for some $v$. For all $y_1, y_2 \in L(g') \cap L(g)$, by Figure 6, we have $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$ and $g\ y_1\ y_2 \Longrightarrow_e y_{12}$ for some $y_{12}$. By Definition B.7, $g' \equiv_\cap g$.

<u>Case 5.8:</u> $g' = \textbf{fuse } d'\ b$ where $d' \neq d$. The proof is similar to the proof of Case 1.5.

<u>Case 5.9:</u> $g' = \textbf{fuse } d\ b$. The proof is similar to the proof of Case 5.7.

<u>Case 6:</u> $g = g_{\text{fa}}$. The proof is by induction on the derivation of $g'$.

<u>Case 6.1:</u> $g' = \textbf{add}$. The proof is similar to the proof of Case 1.5.

<u>Case 6.2:</u> $g' = \textbf{concat}$. The proof is similar to the proof of Case 2.7.

<u>Case 6.3:</u> $g' = \textbf{first}$. The proof is similar to the proof of Case 3.7.

<u>Case 6.4:</u> $g' = \textbf{second}$. The proof is similar to the proof of Case 3.7.

<u>Case 6.5:</u> $g' = \textbf{front } d'\ b$ where $d' \neq d$. The proof is similar to the proof of Case 1.5.

<u>Case 6.6:</u> $g' = \textbf{front } d\ b$. The proof is similar to the proof of Case 5.7.

<u>Case 6.7:</u> $g' = \textbf{back } d'\ b$ where $d' \neq d$. The proof is similar to the proof of Case 1.5.

<u>Case 6.8:</u> $g' = \textbf{back } d\ b$. The proof is similar to the proof of Case 5.7.

<u>Case 6.9:</u> $g' = \textbf{fuse } d'\ b$ where $d' \neq d$. The proof is similar to the proof of Case 1.5.

<u>Case 6.10:</u> $g' = \textbf{fuse } d\ b$. By Definition B.8 and Definition B.1, $y_1, y_2 \in L(g) \cap L(g') = \{y_1' \mathbin{++} d \mathbin{++} y_2' \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} y_k' \mid y_1' \neq \text{nil}, y_k' \neq \text{nil}$, and $y_i' \in L(g_a) \cap L(b)$ and $d \notin y_i'$ for all $i = 1, \ldots, k$, where $k \geq 2\}$. Let $y_1 = y_1^1 \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} y_k^1$ for some $y_i^1 \in L(b)$, where $d \notin y_i^1$ for all $i = 1, \ldots, k$. By Lemma B.3, $y_2 = y_1^2 \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} y_k^2$ and $y_{12} = y_1^{12} \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} y_k^{12}$ for some $y_i^2 \in L(b)$, $y_i^{12} \in \text{String}$, where $d \notin y_i^1$ and $d \notin y_i^2$ for all $i = 1, \ldots, k$. By Definition B.8, $g\ y_1\ y_2 \Longrightarrow_e y_{12}$ and $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $g_a\ y_i^1\ y_i^2 \Longrightarrow_e y_i^{12}$ and $b\ y_i^1\ y_i^2 \Longrightarrow_e y_i^{12}$ for all $i = 1, \ldots, k$.

Let $Y' = \{\langle y_i^1, y_i^2, y_i^{12} \rangle \mid \langle (y_1^1 \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} y_k^1), (y_1^2 \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} y_k^2), (y_1^{12} \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} y_k^{12}) \rangle \in Y$, where $d \notin y_j^1$ and $d \notin y_j^2$ for all $j = 1, \ldots, k$, and $i \in \{1, \ldots, k\}\}$. We have $g_a\ y_i^1\ y_i^2 \Longrightarrow_e y_i^{12}$ and $b\ y_i^1\ y_i^2 \Longrightarrow_e y_i^{12}$ for all $\langle y_i^1, y_i^2, y_i^{12} \rangle \in Y'$. By Definition B.8, $P(g_a, Y')$ and $P(b, Y')$.

By Table 2, $E(g_a, Y')$. By the induction hypothesis, $b \equiv_\cap g_a$.

For all $y_1', y_2' \in L(b) \cap L(g_a)$, by Definition B.7, $b\ y_1'\ y_2' \Longrightarrow_e v$ and $g_a\ y_1'\ y_2' \Longrightarrow_e v$ for some $v$. For all $y_1, y_2 \in L(g') \cap L(g)$, by Figure 6, we have $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$ and $g\ y_1\ y_2 \Longrightarrow_e y_{12}$ for some $y_{12}$. By Definition B.7, $g' \equiv_\cap g$.

<u>Case 7:</u> $g = g_{\text{bfa}}$. The proof is similar to the proof of Case 5.

<u>Case 8:</u> $g = g_{\text{fbfa}}$. The proof is similar to the proof of Case 5.

<u>Case 9:</u> $g = g_{\text{fc}}$. By Table 2, there exists $\langle y_1, y_2, y_{12} \rangle \in Y$. By Definition B.8 and Definition B.1, $y_1, y_2 \in L(g) = \{d \mathbin{++} y \mid y \in \text{String}\}$. Let $y_1 = d \mathbin{++} y_1'$, $y_2 = d \mathbin{++} y_2'$ where $y_1', y_2' \in \text{String}$. By Figure 3, $d \in \text{Delim}$ and $b \in \text{RecOp}$. By Definition B.8, $g\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, there exists $v \in \text{String}$ such that $y_{12} = d \mathbin{++} v$ and $g_c\ y_1'\ y_2' \Longrightarrow_e v$. Hence $v = y_1' \mathbin{++} y_2'$ and $y_{12} = d \mathbin{++} y_1' \mathbin{++} y_2'$.

<u>Case 9.1:</u> $g' = \textbf{add}$. The proof is similar to the proof of Case 1.5.

<u>Case 9.2:</u> $g' = \textbf{concat}$. We show that $P(g', Y)$ never holds in this case. Assume the opposite that $P(g', Y)$. By Definition B.8, $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = y_1 \mathbin{++} y_2 = d \mathbin{++} y_1' \mathbin{++} d \mathbin{++} y_2'$. Since $d \neq \text{nil}$, we have the desired contradiction.

<u>Case 9.3:</u> $g' = \textbf{first}$. We show that $P(g', Y)$ never holds in this case. By Table 2, there exists $y_1', y_2', y_{12}' \in \text{String}$ such that $\langle (d \mathbin{++} y_1'), (d \mathbin{++} y_2'), (d \mathbin{++} y_{12}') \rangle \in Y$ and $y_2' \neq \text{nil}$. By Definition B.8, $g\ (d \mathbin{++} y_1')\ (d \mathbin{++} y_2') \Longrightarrow_e (d \mathbin{++} y_{12}')$. By Figure 6, $y_{12}' = y_1' \mathbin{++} y_2'$. Assume the opposite that $P(g', Y)$. By Definition B.8, $g'\ (d \mathbin{++} y_1')\ (d \mathbin{++} y_2') \Longrightarrow_e (d \mathbin{++} y_{12}')$. By Figure 6, $y_{12}' = y_1'$. Since $y_2' \neq \text{nil}$, we have the desired contradiction.

Case 9.4: $g' = \textbf{second}$. The proof is similar to the proof of Case 9.3.

Case 9.5: $g' = \textbf{front}\ d'\ b$ where $d' \neq d$. The proof is similar to the proof of Case 1.5.

Case 9.6: $g' = \textbf{front}\ d\ b$. The proof is similar to the proof of Case 5.7.

Case 9.7: $g' = \textbf{back}\ d'\ b$ where $d' \neq d$. We show that $P(g', Y)$ never holds in this case. By Definition B.8 and Definition B.1, $y_1, y_2 \in L(g) \cap L(g') = \{d \mathbin{++} y \mid y \in L(g_c)\} \cap \{y \mathbin{++} d' \mid y \in L(b)\}$. For all $\langle y_1, y_2, y_{12} \rangle \in Y$ there exists $y_1', y_2' \in$ String such that $y_1 = d \mathbin{++} y_1' \mathbin{++} d'$, $y_2 = d \mathbin{++} y_2' \mathbin{++} d'$, $(y_1' \mathbin{++} d') \in L(g_c)$, $(y_2' \mathbin{++} d') \in L(g_c)$, $(d \mathbin{++} y_1') \in L(b)$, and $(d \mathbin{++} y_2') \in L(b)$.

By Definition B.8, $g\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = d \mathbin{++} y_1' \mathbin{++} d' \mathbin{++} y_2' \mathbin{++} d'$. By Definition B.8, $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $b\ (d \mathbin{++} y_1')\ (d \mathbin{++} y_2') \Longrightarrow_e (d \mathbin{++} y_1' \mathbin{++} d' \mathbin{++} y_2')$.

Since $d' \neq d$, we have $C(d', (d \mathbin{++} y_1')) = C(d', y_1')$, $C(d', (d \mathbin{++} y_2')) = C(d', y_2')$, and $C(d', (d \mathbin{++} y_1' \mathbin{++} d' \mathbin{++} y_2')) = C(d', y_1') + 1 + C(d', y_2') > C(d', y_1') + C(d', y_2')$, By Figure 3, $d' \in$ Delim and $b \in$ RecOp. By Lemma B.4, we have the desired contradiction.

Case 9.8: $g' = \textbf{back}\ d\ b$. By Definition B.8 and Definition B.1, $y_1, y_2 \in L(g) \cap L(g') = \{d \mathbin{++} y \mid y \in L(g_c)\} \cap \{y \mathbin{++} d \mid y \in L(b)\}$. For all $\langle y_1, y_2, y_{12} \rangle \in Y$ there exists $y_1', y_2' \in$ String such that $y_1 = d \mathbin{++} y_1' \mathbin{++} d$, $y_2 = d \mathbin{++} y_2' \mathbin{++} d$, $(y_1' \mathbin{++} d) \in L(g_c)$, $(y_2' \mathbin{++} d) \in L(g_c)$, $(d \mathbin{++} y_1') \in L(b)$, and $(d \mathbin{++} y_2') \in L(b)$.

By Definition B.8, $g\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = d \mathbin{++} y_1' \mathbin{++} d \mathbin{++} y_2' \mathbin{++} d$. By Definition B.8, $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, $b\ (d \mathbin{++} y_1')\ (d \mathbin{++} y_2') \Longrightarrow_e (d \mathbin{++} y_1' \mathbin{++} d \mathbin{++} y_2')$.

Let $Y' = \{\langle y_1'', y_2'', y_{12}'' \rangle \mid \langle (y_1'' \mathbin{++} d), (y_2'' \mathbin{++} d), (y_{12}'' \mathbin{++} d) \rangle \in Y\}$. We have $b\ y_1''\ y_2'' \Longrightarrow_e y_{12}''$ and $y_{12}'' = y_1'' \mathbin{++} y_2''$ for all $\langle y_1'', y_2'', y_{12}'' \rangle \in Y'$. By Definition B.8, $P(g_c, Y')$ and $P(b, Y')$.

Since $Y \neq \emptyset$, we have $Y' \neq \emptyset$. Also, for all $\langle y_1'', y_2'', y_{12}'' \rangle \in Y'$ there exists $y_1', y_2' \in$ String such that $y_1'' = d \mathbin{++} y_1'$ and $y_2'' = d \mathbin{++} y_2'$. Hence $y_1'' \neq$ nil and $y_2'' \neq$ nil.

By Table 2, $E(g_c, Y')$. By the induction hypothesis, $b \equiv_\cap g_c$.

For all $y_1'', y_2'' \in L(b) \cap L(g_c)$, by Definition B.7, $b\ y_1''\ y_2'' \Longrightarrow_e v$ and $g_c\ y_1''\ y_2'' \Longrightarrow_e v$ for some $v$. By Figure 6, we have $b\ y_1''\ y_2'' \Longrightarrow_e (y_1'' \mathbin{++} y_2'')$ for all $y_1'', y_2'' \in L(b) \cap L(g_c)$. By Definition B.1, $L(g_c) =$ String and $L(b) \cap L(g_c) = L(b)$. Hence $b\ y_1''\ y_2'' \Longrightarrow_e (y_1'' \mathbin{++} y_2'')$ for all $y_1'', y_2'' \in L(b)$.

For all $y_1, y_2 \in L(g) \cap L(g')$, by Definition B.1, $y_1 = d \mathbin{++} y_1' \mathbin{++} d$ and $y_2 = d \mathbin{++} y_2' \mathbin{++} d$ for some $y_1', y_2' \in$ String. Also, $(d \mathbin{++} y_1') \in L(b)$ and $(d \mathbin{++} y_2') \in L(b)$. Hence $b\ (d \mathbin{++} y_1')\ (d \mathbin{++} y_2') \Longrightarrow_e (d \mathbin{++} y_1' \mathbin{++} d \mathbin{++} y_2')$. By Figure 6, $g'\ y_1\ y_2 \Longrightarrow_e d \mathbin{++} y_1' \mathbin{++} d \mathbin{++} y_2' \mathbin{++} d$. Hence $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$ and $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$, where $y_{12} = d \mathbin{++} y_1' \mathbin{++} d \mathbin{++} y_2' \mathbin{++} d$, for all $y_1, y_2 \in L(g) \cap L(g')$. By Definition B.7, $g' \equiv_\cap g$.

Case 9.9: $g' = \textbf{fuse}\ d'\ b$ where $d' \neq d$. We show that $P(g', Y)$ never holds in this case. Recall that $y_1 = d \mathbin{++} y_1'$ and $y_2 = d \mathbin{++} y_2'$ for some $y_1', y_2' \in$ String. Since $d' \neq d$, we have $C(d', y_1) = C(d', y_1')$ and $C(d', y_2) = C(d', y_2')$.

By Definition B.8, $g\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, there exists $v \in$ String such that $y_{12} = d \mathbin{++} v$ and $g_c\ y_1'\ y_2' \Longrightarrow_e v$. We have $v = y_1' \mathbin{++} y_2'$ and $y_{12} = d \mathbin{++} y_1' \mathbin{++} y_2'$. Since $d' \neq d$, we have $C(d', y_{12}) = C(d', y_1') + C(d', y_2') = C(d', y_1) + C(d', y_2)$.

Assume the opposite that $P(g', Y)$. By Definition B.8, we have $g'\ y_1\ y_2 \Longrightarrow_e y_{12}$. By Figure 6, we have $C(d', y_1) = C(d', y_2) \geq 1$. By Figure 3, $d' \in$ Delim and $b \in$ RecOp. By Lemma B.3, $C(d', y_{12}) = C(d', y_1) = C(d', y_2) < C(d', y_1) + C(d', y_2)$. We have the desired contradiction.

Case 9.10: $g' = \textbf{fuse}\ d\ b$. We show that $P(g', Y)$ never holds in this case. Assume the opposite that $P(g', Y)$. By Definition B.8 and Definition B.1, $y_1, y_2 \in L(g') = \{y_1' \mathbin{++} d \mathbin{++} y_2' \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} y_k' \mid y_1' \neq$ nil, $y_k' \neq$ nil, and $y_i' \in L(b)$ and $d \notin y_i'$ for all $i = 1, \ldots, k$, where $k \geq 2\}$. Let $y_1 = y_1^1 \mathbin{++} d \mathbin{++} \ldots \mathbin{++} d \mathbin{++} y_k^1$ for some $y_1^1, \ldots, y_k^1 \in L(b)$. We have $y_1^1 \neq$ nil and $d \notin y_1^1$. Recall that $y_1 = d \mathbin{++} y_1'$. We have the desired contradiction.

$\square$

**Theorem 7.** For any combiners $g \in G_{\text{rec}}, g' \in$ RecOp and set of output tuples $Y$, if $E_{\text{rec}}(Y)$, $P(g, Y)$, and $P(g', Y)$, then $g' \equiv_\cap g$.

*Proof.* By Theorem 6 and Proposition B.2. $\square$

**Theorem 8.** For any command $f$, set of input streams $X$, combiner $g \in G_{\text{rec}}$, and combiner $g' \in$ RecOp, if the following conditions hold:

- $E_{\text{rec}}(f(X))$,
- $g$ is correct for $f$, and

- $y_1, y_2 \in L(g')$ for all $\langle y_1, y_2, y_{12} \rangle \in f(X)$,

then $P(g', f(X))$ if and only if $g' \equiv_\cap g$.

*Proof.* To show that "if $P(g', f(X))$ then $g' \equiv_\cap g$": By Definition B.9, $P(g, f(X))$. By Theorem 7, $g' \equiv_\cap g$.

To show that "if $g' \equiv_\cap g$ then $P(g', f(X))$": For all $\langle y_1, y_2, y_{12} \rangle \in f(X)$, by Proposition B.1, $y_1, y_2 \in L(g)$. Hence $y_1, y_2 \in L(g) \cap L(g')$. By Definition B.7, $g\, y_1\, y_2 \Longrightarrow_e v$ and $g'\, y_1\, y_2 \Longrightarrow_e v$ for some $v$. By Definition B.9, $P(g, f(X))$. By Definition B.8, $g\, y_1\, y_2 \Longrightarrow_e y_{12}$. Hence $v = y_{12}$. By Definition B.8, $P(g', f(X))$. □

**Theorem 9.** For any command $f$, set of input streams $X$, combiner $g \in G_{\text{rec}}$, combiner $g' \in \text{RecOp}$, and integer $k$, if the following conditions hold:

- $E_{\text{rec}}(f(X))$,
- $g$ is correct for $f$,
- $y_1, y_2 \in L(g')$ for all $\langle y_1, y_2, y_{12} \rangle \in f(X)$, and
- $k \geq |g'|$,

then $g' \in P_k(f(X)) \cap \text{RecOp}$ if and only if $g' \equiv_\cap g$.

*Proof.* By Theorem 8 and Definition B.17. □

*Remark.* Theorem 9 states that, if the specified size is large enough, if correct combiner is among $G_{\text{rec}}$, and if the synthesizer has collected sufficient observations, then the synthesizer must return either the correct combiner or its equivalent. By Proposition B.7, we know that as long as the specified size $k \geq 6$, the synthesizer is guaranteed to return a correct combiner if the correct combiner is among $G_{\text{rec}}$.

**Lemma B.5.** For any set of output tuples $Y$, let $E(Y)$ be a predicate that is true if and only if the following conditions hold:

- There exists $\langle y_1, y_2, y_{12} \rangle \in Y$ and $c \in y_1$ such that $c \notin \text{Delim} \cup \{\text{`0`}\}$.
- There exists $\langle y_1, y_2, y_{12} \rangle \in Y$ and $c \in y_2$ such that $c \notin \text{Delim} \cup \{\text{`0`}\}$.

For any combiner $g \in \{g_f, g_s\}$, set of output tuples $Y$ such that $P(g, Y)$ and $E(Y)$, and $g' \in \text{RecOp}$, if $P(g', Y)$ then either $g' \equiv_\cap g_f$ or $g' \equiv_\cap g_s$.

*Proof.* <u>Case 1:</u> There exists $\langle y_1, y_2, y_{12} \rangle \in Y$ such that $y_1 \neq y_2$. By Theorem 6, $g' \equiv_\cap g$.
<u>Case 2:</u> For all $\langle y_1, y_2, y_{12} \rangle \in Y$, $y_1 = y_2$. The proof is similar to the proof of Case 3 in Theorem 6. □

**Theorem 10.** For any combiner $g \in G_{\text{struct}}$, set of output tuples $Y$ such that $P(g, Y)$ and $E(g, Y)$, and $g' \in \text{StructOp}$, we have $P(g', Y)$ implies $g' \equiv_\cap g$.

*Proof.* The proof performs case analysis of the values of $g, g'$.

<u>Case 1:</u> $g = g_{\text{sf}}$: By Table 2, there exists $\langle y_1, y_2, y_{12} \rangle \in Y$ such that $(\text{splitLastLine } y_1) = (y_1', l)$ and $(\text{splitFirstLine } y_2) = (l, y_2')$ for some $l$, where $(\text{firstChar } (\text{delPad } l)) \notin \text{Delim} \cup \{\text{`0`}\}$ and $(\text{lastChar } l) \notin \text{Delim} \cup \{\text{`0`}\}$.

    <u>Case 1.1:</u> $g' = \textbf{stitch } b$: By Definition B.8, $g\, y_1\, y_2 \Longrightarrow_e y_{12}$ and $g'\, y_1\, y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = y_1'\, {+\!\!+}$ `\n` ${+\!\!+}\, v\, {+\!\!+}$ `\n` ${+\!\!+}\, y_2'$ for some $v$, where $g_f\, l\, l \Longrightarrow_e v$ and $b\, l\, l \Longrightarrow_e v$. By Figure 6, $v = l$ and $b\, l\, l \Longrightarrow_e l$.

        Let $c = \text{lastChar } l$. We have $c \in l$ and $c \notin \text{Delim} \cup \{\text{`0`}\}$. Let $Y' = \{\langle l, l, l \rangle\}$. By Lemma B.5, either $b \equiv_\cap g_f$ or $b \equiv_\cap g_s$. Either case, by Definition B.7, we have $b\, l'\, l' \Longrightarrow_e l'$ for all $l' \in L(b) \cap L(g_s) = L(b)$. By Definition B.1, $L(g') = L(g') \cap L(g)$. By Figure 6, for all $y_1, y_2 \in L(g')$ we have $g'\, y_1\, y_2 \Longrightarrow_e v'$ and $g\, y_1\, y_2 \Longrightarrow_e v$ for some $v$. By Definition B.7, $g' \equiv_\cap g$.

    <u>Case 1.2:</u> $g' = \textbf{stitch2 } d\, b_1\, b_2$: By Definition B.8, $g\, y_1\, y_2 \Longrightarrow_e y_{12}$ and $g'\, y_1\, y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = y_1'\, {+\!\!+}$ `\n` ${+\!\!+}\, v\, {+\!\!+}$ `\n` ${+\!\!+}\, y_2'$ for some $v$, where $v = \text{addPad } (h\, {+\!\!+}\, d\, {+\!\!+}\, t)$. Let $(h', t') = \text{splitFirst } d\, (\text{delPad } l)$. We have $d \notin h'$.

        By Figure 6, we have $g_f\, (h'\, {+\!\!+}\, d\, {+\!\!+}\, t')\, (h'\, {+\!\!+}\, d\, {+\!\!+}\, t') \Longrightarrow_e (h\, {+\!\!+}\, d\, {+\!\!+}\, t)$, $b_1\, h'\, h' \Longrightarrow_e h$, and $b_2\, t'\, t' \Longrightarrow_e t$. Hence $(h'\, {+\!\!+}\, d\, {+\!\!+}\, t' = h\, {+\!\!+}\, d\, {+\!\!+}\, t$. By Figure 3, $d \in \text{Delim}$ and $b_1 \in \text{RecOp}$. By Lemma B.1, $d \notin h$. Hence $h' = h$ and $t' = t$. We have $b_1\, h\, h \Longrightarrow_e h$ and $b_2\, t\, t \Longrightarrow_e t$.

        Let $c_1 = (\text{firstChar } (\text{delPad } l))$ and $c_2 = \text{lastChar } l$. We have $c_1 \in h$, $c_2 \in t$, and $c_1, c_2 \notin \text{Delim} \cup \{\text{`0`}\}$.

        By Definition B.8, for all $\langle y_1, y_2, y_{12} \rangle \in Y$ we have $y_1, y_2 \in L(g')$. By Table 2, then there exists $\langle y_1, y_2, y_{12} \rangle \in Y$ such that $h_1 \neq h_2$, where $(\text{splitLastLine } y_1) = (y_1', l_1)$, $(\text{splitFirstLine } y_2) = (l_2, y_2')$, $(\text{splitFirst } d\, (\text{delPad } l_1)) =$

$(h_1, t)$, and (splitFirst $d$ (delPad $l_2$)) $= (h_2, t)$. By Definition B.8, Figure 6, Figure 3, and Lemma B.1, $b_1 \; h_1 \; h_2 \Longrightarrow_e h_1$.

Let $Y_1' = \{\langle h, h, h\rangle, \langle h_1, h_2, h_1\rangle\}$. By Table 2, we have $E(g_f, Y_1')$. By Figure 6, $g_f \; h \; h \Longrightarrow_e h$ and $g_f \; h_1 \; h_2 \Longrightarrow_e h_1$. By Definition B.8, $P(g_f, Y_1')$ and $P(b_1, Y_1')$. By Figure 3, $b_1 \in$ RecOp. By Theorem 6, $b_1 \equiv_\cap g_f$.

Let $Y_2' = \{\langle t, t, t\rangle\}$. By Lemma B.5, either $b_2 \equiv_\cap g_f$ or $b_2 \equiv_\cap g_s$. The rest of the proof is similar to the proof of Case 1.1.

Case 1.3: $g' = $ **offset** $d \; b$: We show that $P(g', Y)$ never holds in this case. Let $n_1, n_2, n_{12}$ be the numbers of lines in $y_1, y_2, y_{12}$, respectively. By Definition B.8, $g \; y_1 \; y_2 \Longrightarrow_e y_{12}$. By Figure 6, $n_{12} = n_1 + n_2 - 1$. Assume the opposite that $P(g', Y)$. By Definition B.8, $g' \; y_1 \; y_2 \Longrightarrow_e y_{12}$. Since $l \neq$ nil, by Figure 6, $n_{12} = n_1 + n_2$. We have the desired contradiction.

Case 2: $g = g_{saf}$: By Table 2, there exists $\langle y_1, y_2, y_{12}\rangle \in Y$ such that (splitLastLine $y_1$) $= (y_1', l)$ and (splitFirstLine $y_2$) $= (l, y_2')$ for some $l$, where (firstChar (delPad $l$)) $\notin$ Delim $\cup \{`0`\}$ and (lastChar $l$) $\notin$ Delim $\cup \{`0`\}$. By Definition B.8 and Definition B.1, $l = p ++ h ++ d ++ t$ for some $p \in [`\;`^{+} \;|\; `\backslash t`]$, $h \in L(g_a) = [`0` - `9`]^{+}$, and $t \in L(g_f) =$ String.

Case 2.1: $g' = $ **stitch** $b$: We outline the proof below. By Definition B.8, $g \; y_1 \; y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = y_1' ++ `\backslash n` ++ v ++ `\backslash n` ++ y_2'$ for some $v$, where $v = p ++ h_{12} ++ d ++ t_{12}$, $g_a \; h \; h \Longrightarrow_e h_{12}$, and $g_f \; t \; t \Longrightarrow_e t_{12}$. We have $h_{12} = $ (intToStr ((strToInt $h$) + (strToInt $h$))) and $t_{12} = t$. By Definition B.8, $g' \; y_1 \; y_2 \Longrightarrow_e y_{12}$. By Figure 6, $b \; l \; l \Longrightarrow_e v$. Such $b \in$ RecOp does not exist. The proof is by induction on the derivation of $b$.

Case 2.2: $g' = $ **stitch2** $d' \; b_1 \; b_2$, where $d' \neq d$: We show that such $b_1$ does not exist. The proof is similar to the proof of Case 2.1.

Case 2.3: $g' = $ **stitch2** $d \; b_1 \; b_2$: We outline the proof below. By Theorem 6, $b_1 \equiv_\cap g_a$. By Lemma B.5, either $b_2 \equiv_\cap g_f$ or $b_2 \equiv_\cap g_s$. Either case, $g' \equiv_\cap g$.

Case 2.4: $g' = $ **offset** $d' \; b$: The proof is similar to the proof of Case 1.3.

Case 3: $g = g_{oa}$: By Table 2, there exists $\langle y_1, y_2, y_{12}\rangle \in Y$ such that (splitLastLine $y_1$) $= (y_1', l_1)$, (splitFirstLine $y_2$) $= (l_2, y_2')$, and (splitFirstLine $y_2'$) $= (l_2', y_2'')$, where (firstChar (delPad $l_1$)) $\notin$ Delim $\cup \{`0`\}$, $l_2 \neq$ nil, and $l_2' \neq$ nil.

Case 3.1: $g' = $ **stitch** $b$: We outline the proof below. We show that $P(g', Y)$ never holds in this case. Assume the opposite that $P(g', Y)$. By Definition B.8, $g' \; y_1 \; y_2 \Longrightarrow_e y_{12}$. If $l_1 = l_2$, the proof is similar to the proof of Case 1.3. If $l_1 \neq l_2$, by Figure 6, $y_{12} = y_1 ++ y_2$. By Definition B.8, $g \; y_1 \; y_2 \Longrightarrow_e y_{12}$. By Figure 6, $y_{12} = y_1 ++ v$ for some $v$. By Table 2, $v \neq y_2$.

Case 3.2: $g' = $ **stitch2** $d' \; b_1 \; b_2$: The proof is similar to the proof of Case 3.1.

Case 3.3: $g' = $ **offset** $d' \; b$, where $d' \neq d$: We show that such $b$ does not exist. The proof is similar to the proof of Case 2.1.

Case 3.4: $g' = $ **offset** $d \; b$: We outline the proof below. Let $Y' = \{\langle h_1, h_2, y_{12}'\rangle \mid \langle y_1, y_2, y_{12}\rangle \in Y,$ (splitLastLine $y_1$) $= (y_1', l_1),$ (splitFirstLine $y_2$) $= (l_2, y_2'),$ (splitFirst $d$ (delPad $l_1$)) $= (h_1, t_1),$ and (splitFirst $d$ (delPad $l_2$)) $= (h_2, t_2)\}$. By Table 2, $E(g_a, Y')$. By Theorem 6, $b \equiv_\cap g_a$.

□

**Theorem 11.** For any combiners $g \in G_{struct}, g' \in$ StructOp and set of output tuples $Y$, if $E_{struct}(Y)$, $P(g, Y)$, and $P(g', Y)$, then $g' \equiv_\cap g$.

*Proof.* By Theorem 10 and Proposition B.4. □

**Theorem 12.** For any command $f$, set of input streams $X$, combiner $g \in G_{struct}$, and combiner $g' \in$ StructOp, if the following conditions hold:

- $E_{struct}(f(X))$,
- $g$ is correct for $f$, and
- $y_1, y_2 \in L(g')$ for all $\langle y_1, y_2, y_{12}\rangle \in f(X)$,

then $P(g', f(X))$ if and only if $g' \equiv_\cap g$.

*Proof.* The proof is similar to the proof of Theorem 8. □

**Theorem 13.** For any command $f$, set of input streams $X$, combiner $g \in G_{\text{struct}}$, combiner $g' \in \text{StructOp}$, and integer $k$, if the following conditions hold:

- $E_{\text{struct}}(f(X))$,
- $g$ is correct for $f$,
- $y_1, y_2 \in L(g')$ for all $\langle y_1, y_2, y_{12} \rangle \in f(X)$, and
- $k \geq |g'|$,

then $g' \in P_k(f(X)) \cap \text{StructOp}$ if and only if $g' \equiv_{\cap} g$.

*Proof.* By Theorem 12 and Definition B.17. □

*Remark.* Theorem 13 states that, if the specified size is large enough, if correct combiner is among $G_{\text{struct}}$, and if the synthesizer has collected sufficient observations and eliminated RecOp candidates, then the synthesizer must return either the correct combiner or its equivalent. By Proposition B.7, we know that as long as the specified size $k \geq 6$, the synthesizer is guaranteed to return a correct combiner if the correct combiner is among $G_{\text{struct}}$.

## C   Appendix: Performance Results

Table 3 presents the pipeline stages that are automatically parallelized by KUMQUAT. The first two columns present the benchmark and script names. The next column (**Parallelized**) presents the number of stages automatically parallelized by KUMQUAT, $k$, and the number of stages in the original pipeline, $n$, as a pair "$k/n$" for each pipeline in the parentheses. Here we also report the single commands that are not in pipelines, as pairs "$k/1$". The pair before the parentheses presents the sum over all pipelines in the script. The next column (**Eliminated**) presents the number of parallelized stages whose combiners are eliminated by KUMQUAT during optimization. Again, the numbers in the parentheses correspond to pipelines in the script. The number before the parentheses presents the sum over all pipelines.

Table 4 compares the parallel execution times with the original script execution times. The first two columns present the benchmark and script names. The next column ($T_{\text{orig}}$) presents the execution time for the original unmodified benchmark script. The next column ($u_1$) presents the serial execution time. The next column ($u_{16}$) presents the *unoptimized* parallel execution time with 16 way parallelism. The next column ($T_{16}$) presents the *optimized* parallel execution time with 16 way parallelism.

Table 5 presents the parallel execution times for *unoptimized* pipelines with 1, 2, 4, 8, and 16 way parallelism.

Table 6 presents the parallel execution times for *optimized* pipelines with 1, 2, 4, 8, and 16 way parallelism.

Among all benchmark scripts, the unoptimized parallel speedup ranges between 0.5× and 14.9×, with a median speedup of 5.3×. The optimized parallel speedup ranges between 0.6× and 26.9×, with a median speedup of 7.1× (we attribute the superlinear speedup to pipelined parallelism exploited across consecutive parallelized commands with no intermediate combiner). All scripts that exhibit a slowdown have a serial execution time under 10 seconds.

Table 7 presents the performance results for benchmark scripts whose serial execution time is at least 3 minutes. In general, shorter scripts have smaller parallel speedup.

**Table 3.** Pipeline commands that are parallelized with the synthesized combiners

| Benchmark | Script Name | Parallelized | Eliminated |
|---|---|---|---|
| analytics-mts | 1.sh (vehicles per day) | 7/7 (7/7) | 3 (3) |
| analytics-mts | 2.sh (vehicle days on road) | 8/8 (8/8) | 3 (3) |
| analytics-mts | 3.sh (vehicle hours on road) | 8/8 (8/8) | 3 (3) |
| analytics-mts | 4.sh (hours monitored per day) | 7/7 (7/7) | 3 (3) |
| oneliners | bi-grams.sh | 3/5 (3/5) | 0 (0) |
| oneliners | diff.sh | 4/7 (0/1, 2/2, 2/2, 0/1, 0/1) | 2 (0, 1, 1, 0, 0) |
| oneliners | nfa-regex.sh | 2/2 (2/2) | 1 (1) |
| oneliners | set-diff.sh | 5/8 (0/1, 3/3, 2/2, 0/1, 0/1) | 3 (0, 2, 1, 0, 0) |
| oneliners | shortest-scripts.sh | 6/7 (6/7) | 5 (5) |
| oneliners | sort-sort.sh | 3/3 (3/3) | 1 (1) |
| oneliners | sort.sh | 1/1 (1/1) | 0 (0) |
| oneliners | spell.sh | 6/8 (6/8) | 3 (3) |
| oneliners | top-n.sh | 4/6 (4/6) | 1 (1) |
| oneliners | wf.sh | 4/5 (4/5) | 1 (1) |
| poets | 1_1.sh (count_words) | 4/6 (4/6) | 1 (1) |
| poets | 2_1.sh (merge_upper) | 5/7 (5/7) | 2 (2) |
| poets | 2_2.sh (count_vowel_seq) | 5/7 (5/7) | 2 (2) |
| poets | 3_1.sh (sort) | 5/7 (5/7) | 1 (1) |

| | | | |
|---|---|---|---|
| poets | 3_2.sh (sort_words_by_folding) | 5/7 (5/7) | 1 (1) |
| poets | 3_3.sh (sort_words_by_rhyming) | 7/9 (7/9) | 2 (2) |
| poets | 4_3.sh (bigrams) | 4/8 (2/4, 0/1, 2/3) | 1 (1, 0, 0) |
| poets | 4_3b.sh (count_trigrams) | 4/9 (2/4, 0/1, 0/1, 2/3) | 1 (1, 0, 0, 0) |
| poets | 6_1.sh (trigram_rec) | 8/14 (4/7, 4/7) | 4 (2, 2) |
| poets | 6_1_1.sh (uppercase_by_token) | 3/5 (3/5) | 1 (1) |
| poets | 6_1_2.sh (uppercase_by_type) | 4/6 (4/6) | 1 (1) |
| poets | 6_2.sh (4letter_words) | 7/11 (3/5, 4/6) | 2 (1, 1) |
| poets | 6_3.sh (words_no_vowels) | 5/7 (5/7) | 2 (2) |
| poets | 6_4.sh (1syllable_words) | 5/8 (5/8) | 2 (2) |
| poets | 6_5.sh (2syllable_words) | 5/8 (5/8) | 2 (2) |
| poets | 6_7.sh (verses_2om_3om_2instances) | 10/13 (3/4, 3/4, 4/5) | 7 (2, 2, 3) |
| poets | 7_2.sh (count_consonant_seq) | 5/7 (5/7) | 2 (2) |
| poets | 8.2_1.sh (vowel_sequences_gr_1K) | 5/8 (5/8) | 1 (1) |
| poets | 8.2_2.sh (bigrams_appear_twice) | 4/9 (2/4, 0/1, 2/3, 0/1) | 1 (1, 0, 0, 0) |
| poets | 8.3_2.sh (find_anagrams) | 7/9 (2/4, 1/1, 1/1, 3/3) | 1 (1, 0, 0, 0) |
| poets | 8.3_3.sh (compare_exodus_genesis) | 6/10 (3/5, 1/2, 2/3) | 1 (1, 0, 0) |
| poets | 8_1.sh (sort_words_by_n_syllables) | 6/10 (3/5, 2/2, 1/3) | 2 (1, 1, 0) |
| unix50 | 1.sh (1.0: extract last name) | 1/1 (1/1) | 0 (0) |
| unix50 | 10.sh (4.4: histogram by piece) | 9/9 (9/9) | 6 (6) |
| unix50 | 11.sh (4.5: histogram by piece and pawn) | 9/9 (9/9) | 6 (6) |
| unix50 | 12.sh (4.6: piece used most) | 8/9 (8/9) | 5 (5) |
| unix50 | 13.sh (5.1: extract hellow world) | 3/3 (3/3) | 2 (2) |
| unix50 | 14.sh (6.1: order bodies) | 3/3 (3/3) | 1 (1) |
| unix50 | 15.sh (7.1: number of versions) | 3/3 (3/3) | 2 (2) |
| unix50 | 16.sh (7.2: most frequent machine) | 6/7 (6/7) | 1 (1) |
| unix50 | 17.sh (7.3: decades unix released) | 5/5 (5/5) | 2 (2) |
| unix50 | 18.sh (8.1: count unix birth-year) | 3/3 (3/3) | 2 (2) |
| unix50 | 19.sh (8.2: location office) | 4/4 (4/4) | 3 (3) |
| unix50 | 2.sh (1.1: extract names and sort) | 2/2 (2/2) | 1 (1) |
| unix50 | 20.sh (8.3: four most involved) | 4/4 (4/4) | 3 (3) |
| unix50 | 21.sh (8.4: longest words w/o hyphens) | 3/3 (3/3) | 1 (1) |
| unix50 | 23.sh (9.1: extract word PORT) | 6/6 (6/6) | 4 (4) |
| unix50 | 24.sh (9.2: extract word BELL) | 2/2 (2/2) | 1 (1) |
| unix50 | 25.sh (9.3: animal decorate) | 2/2 (2/2) | 1 (1) |
| unix50 | 26.sh (9.4: four corners) | 4/5 (4/5) | 2 (2) |
| unix50 | 28.sh (9.6: follow directions) | 6/10 (6/10) | 3 (3) |
| unix50 | 29.sh (9.7: four corners) | 2/4 (2/4) | 1 (1) |
| unix50 | 3.sh (1.2: extract names and sort) | 1/2 (1/2) | 0 (0) |
| unix50 | 30.sh (9.8: TELE-communications) | 4/8 (4/8) | 2 (2) |
| unix50 | 31.sh (9.9) | 4/9 (4/9) | 2 (2) |
| unix50 | 32.sh (10.1: count recipients) | 3/4 (3/4) | 2 (2) |
| unix50 | 33.sh (10.2: list recipients) | 2/3 (2/3) | 1 (1) |
| unix50 | 34.sh (10.3: extract username) | 7/7 (7/7) | 4 (4) |
| unix50 | 35.sh (11.1: year received medal) | 2/2 (2/2) | 1 (1) |
| unix50 | 36.sh (11.2: most repeated first name) | 7/8 (7/8) | 2 (2) |
| unix50 | 4.sh (1.3: sort top first names) | 4/4 (4/4) | 1 (1) |
| unix50 | 5.sh (2.1: all Unix utilities) | 2/2 (2/2) | 1 (1) |
| unix50 | 6.sh (3.1: first letter of last names) | 4/4 (4/4) | 2 (2) |
| unix50 | 7.sh (4.1: number of rounds) | 3/3 (3/3) | 2 (2) |
| unix50 | 8.sh (4.2: pieces captured) | 4/4 (4/4) | 3 (3) |
| unix50 | 9.sh (4.3: pieces captured with pawn) | 6/6 (6/6) | 5 (5) |
| **Total** | | 325/427 | 144 |

**Table 4.** Performance results for all benchmark scripts, comparing new pipelines with original scripts

| Benchmark | Script Name | $T_{\mathbf{orig}}$ | $u_1$ | $u_{16}$ | $T_{16}$ |
|---|---|---|---|---|---|
| analytics-mts | 1.sh (vehicles per day) | 333 s (1.1×) | 376 s | 40 s (9.4×) | 29 s (13.1×) |
| analytics-mts | 2.sh (vehicle days on road) | 335 s (1.1×) | 379 s | 41 s (9.3×) | 28 s (13.5×) |
| analytics-mts | 3.sh (vehicle hours on road) | 408 s (1.0×) | 427 s | 51 s (8.4×) | 38 s (11.3×) |
| analytics-mts | 4.sh (hours monitored per day) | 99 s (1.7×) | 167 s | 28 s (6.0×) | 13 s (12.8×) |
| oneliners | bi-grams.sh | 668 s (1.5×) | 1007 s | 118 s (8.6×) | 115 s (8.7×) |

| | | | | | |
|---|---|---|---|---|---|
| oneliners | diff.sh | 325 s (1.5×) | 478 s | 98 s (4.9×) | 83 s (5.8×) |
| oneliners | nfa-regex.sh | 389 s (1.0×) | 391 s | 26 s (14.9×) | 27 s (14.7×) |
| oneliners | set-diff.sh | 879 s (1.5×) | 1308 s | 144 s (9.1×) | 128 s (10.2×) |
| oneliners | shortest-scripts.sh | 82 s (1.3×) | 110 s | 9 s (12.4×) | 7 s (16.2×) |
| oneliners | sort-sort.sh | 137 s (1.2×) | 167 s | 31 s (5.4×) | 28 s (6.0×) |
| oneliners | sort.sh | 273 s (1.4×) | 389 s | 39 s (10.0×) | 38 s (10.3×) |
| oneliners | spell.sh | 427 s (1.7×) | 736 s | 78 s (9.5×) | 61 s (12.1×) |
| oneliners | top-n.sh | 372 s (1.7×) | 622 s | 63 s (9.9×) | 50 s (12.4×) |
| oneliners | wf.sh | 1155 s (1.8×) | 2089 s | 196 s (10.7×) | 145 s (14.4×) |
| poets | 1_1.sh (count_words) | 360 s (1.8×) | 637 s | 84 s (7.6×) | 83 s (7.6×) |
| poets | 2_1.sh (merge_upper) | 307 s (1.8×) | 547 s | 79 s (6.9×) | 78 s (7.0×) |
| poets | 2_2.sh (count_vowel_seq) | 112 s (1.2×) | 140 s | 27 s (5.2×) | 24 s (5.8×) |
| poets | 3_1.sh (sort) | 391 s (1.7×) | 665 s | 89 s (7.4×) | 88 s (7.6×) |
| poets | 3_2.sh (sort_words_by_folding) | 402 s (1.7×) | 681 s | 94 s (7.3×) | 94 s (7.2×) |
| poets | 3_3.sh (sort_words_by_rhyming) | 415 s (1.7×) | 699 s | 100 s (7.0×) | 100 s (7.0×) |
| poets | 4_3.sh (bigrams) | 635 s (1.4×) | 915 s | 173 s (5.3×) | 173 s (5.3×) |
| poets | 4_3b.sh (count_trigrams) | 862 s (1.2×) | 1049 s | 275 s (3.8×) | 279 s (3.8×) |
| poets | 6_1.sh (trigram_rec) | 2 s (1.9×) | 5 s | 8 s (0.6×) | 2 s (2.4×) |
| poets | 6_1_1.sh (uppercase_by_token) | 38 s (1.2×) | 45 s | 14 s (3.3×) | 14 s (3.2×) |
| poets | 6_1_2.sh (uppercase_by_type) | 330 s (1.9×) | 635 s | 64 s (10.0×) | 24 s (26.9×) |
| poets | 6_2.sh (4letter_words) | 327 s (2.0×) | 647 s | 80 s (8.1×) | 34 s (18.8×) |
| poets | 6_3.sh (words_no_vowels) | 220 s (1.1×) | 235 s | 32 s (7.4×) | 31 s (7.7×) |
| poets | 6_4.sh (1syllable_words) | 433 s (1.3×) | 542 s | 57 s (9.5×) | 31 s (17.4×) |
| poets | 6_5.sh (2syllable_words) | 397 s (1.1×) | 443 s | 48 s (9.2×) | 40 s (11.0×) |
| poets | 6_7.sh (verses_2om_3om_2instances) | 4 s (2.0×) | 7 s | 11 s (0.6×) | 5 s (1.5×) |
| poets | 7_2.sh (count_consonant_seq) | 475 s (1.4×) | 678 s | 80 s (8.5×) | 48 s (14.2×) |
| poets | 8.2_1.sh (vowel_sequences_gr_1K) | 417 s (1.4×) | 573 s | 73 s (7.9×) | 42 s (13.7×) |
| poets | 8.2_2.sh (bigrams_appear_twice) | 645 s (1.4×) | 921 s | 177 s (5.2×) | 91 s (10.2×) |
| poets | 8.3_2.sh (find_anagrams) | 237 s (3.1×) | 724 s | 102 s (7.1×) | 50 s (14.5×) |
| poets | 8.3_3.sh (compare_exodus_genesis) | 334 s (2.0×) | 656 s | 74 s (8.8×) | 34 s (19.3×) |
| poets | 8_1.sh (sort_words_by_n_syllables) | 346 s (1.9×) | 653 s | 69 s (9.5×) | 26 s (24.6×) |
| unix50 | 1.sh (1.0: extract last name) | 12 s (1.0×) | 12 s | 3 s (3.6×) | 3 s (3.5×) |
| unix50 | 10.sh (4.4: histogram by piece) | 27 s (1.8×) | 48 s | 13 s (3.7×) | 6 s (7.8×) |
| unix50 | 11.sh (4.5: histogram by piece and pawn) | 25 s (1.7×) | 42 s | 12 s (3.4×) | 6 s (6.8×) |
| unix50 | 12.sh (4.6: piece used most) | 115 s (1.3×) | 149 s | 27 s (5.5×) | 18 s (8.2×) |
| unix50 | 13.sh (5.1: extract hellow world) | 4 s (2.7×) | 12 s | 7 s (1.7×) | 2 s (5.1×) |
| unix50 | 14.sh (6.1: order bodies) | 143 s (1.3×) | 185 s | 31 s (6.0×) | 25 s (7.5×) |
| unix50 | 15.sh (7.1: number of versions) | 5 s (1.4×) | 8 s | 6 s (1.4×) | 3 s (2.5×) |
| unix50 | 16.sh (7.2: most frequent machine) | 80 s (1.2×) | 93 s | 15 s (6.4×) | 13 s (7.4×) |
| unix50 | 17.sh (7.3: decades unix released) | 39 s (1.1×) | 43 s | 10 s (4.4×) | 8 s (5.1×) |
| unix50 | 18.sh (8.1: count unix birth-year) | 2 s (1.7×) | 3 s | 6 s (0.5×) | 3 s (1.1×) |
| unix50 | 19.sh (8.2: location office) | 2 s (1.3×) | 2 s | 2 s (1.0×) | 2 s (1.0×) |
| unix50 | 2.sh (1.1: extract names and sort) | 133 s (1.3×) | 171 s | 20 s (8.4×) | 17 s (9.9×) |
| unix50 | 20.sh (8.3: four most involved) | 0 s (NaN) | 5 s | 6 s (0.8×) | 3 s (1.7×) |
| unix50 | 21.sh (8.4: longest words w/o hyphens) | 428 s (1.7×) | 733 s | 64 s (11.4×) | 49 s (14.9×) |
| unix50 | 23.sh (9.1: extract word PORT) | 111 s (1.8×) | 202 s | 23 s (8.8×) | 10 s (19.8×) |
| unix50 | 24.sh (9.2: extract word BELL) | 4 s (1.1×) | 5 s | 2 s (2.1×) | 2 s (2.4×) |
| unix50 | 25.sh (9.3: animal decorate) | 5 s (1.1×) | 6 s | 3 s (2.1×) | 2 s (2.3×) |
| unix50 | 26.sh (9.4: four corners) | 11 s (2.8×) | 32 s | 18 s (1.7×) | 15 s (2.1×) |
| unix50 | 28.sh (9.6: follow directions) | 87 s (2.2×) | 188 s | 54 s (3.5×) | 49 s (3.8×) |
| unix50 | 29.sh (9.7: four corners) | 6 s (3.0×) | 19 s | 18 s (1.0×) | 15 s (1.3×) |
| unix50 | 3.sh (1.2: extract names and sort) | 0 s (NaN) | 0 s | 0 s (0.7×) | 0 s (0.7×) |
| unix50 | 30.sh (9.8: TELE-communications) | 100 s (1.6×) | 154 s | 66 s (2.3×) | 62 s (2.5×) |
| unix50 | 31.sh (9.9) | 88 s (1.7×) | 149 s | 73 s (2.0×) | 68 s (2.2×) |
| unix50 | 32.sh (10.1: count recipients) | 3 s (1.8×) | 6 s | 7 s (0.9×) | 6 s (0.9×) |
| unix50 | 33.sh (10.2: list recipients) | 3 s (1.8×) | 6 s | 6 s (1.0×) | 5 s (1.1×) |
| unix50 | 34.sh (10.3: extract username) | 0 s (NaN) | 2 s | 3 s (0.7×) | 3 s (0.9×) |
| unix50 | 35.sh (11.1: year received medal) | 1 s (0.9×) | 1 s | 2 s (0.5×) | 2 s (0.6×) |
| unix50 | 36.sh (11.2: most repeated first name) | 15 s (1.3×) | 19 s | 7 s (2.8×) | 6 s (3.5×) |
| unix50 | 4.sh (1.3: sort top first names) | 134 s (1.1×) | 154 s | 21 s (7.4×) | 19 s (8.3×) |
| unix50 | 5.sh (2.1: all Unix utilities) | 7 s (1.1×) | 8 s | 3 s (2.5×) | 2 s (3.1×) |
| unix50 | 6.sh (3.1: first letter of last names) | 10 s (1.4×) | 14 s | 5 s (2.7×) | 3 s (5.2×) |
| unix50 | 7.sh (4.1: number of rounds) | 15 s (1.2×) | 18 s | 9 s (2.0×) | 4 s (4.7×) |
| unix50 | 8.sh (4.2: pieces captured) | 6 s (2.1×) | 12 s | 8 s (1.6×) | 3 s (4.1×) |

| | | | | | | |
|---|---|---|---|---|---|---|
| unix50 | 9.sh (4.3: pieces captured with pawn) | | 14 s (2.0×) | 28 s | 9 s (2.9×) | 4 s (7.3×) |
| **Max** | | | 1155 s (3.1×) | 2089 s | 275 s (14.9×) | 279 s (26.9×) |
| **Min** | | | 0 s (0.9×) | 0 s | 0 s (0.5×) | 0 s (0.6×) |
| **Mean** | | | 217 s (1.6×) | 332 s | 48 s (5.5×) | 37 s (8.0×) |
| **Median** | | | 114 s (1.5×) | 167 s | 28 s (5.3×) | 24 s (7.1×) |

**Table 5.** Performance results for all benchmark scripts, where new pipelines are unoptimized

| Benchmark | Script Name | $u_1$ | $u_2$ | $u_4$ | $u_8$ | $u_{16}$ |
|---|---|---|---|---|---|---|
| analytics-mts | 1.sh (vehicles per day) | 376 s | 200 s (1.9×) | 107 s (3.5×) | 62 s (6.1×) | 40 s (9.4×) |
| analytics-mts | 2.sh (vehicle days on road) | 379 s | 199 s (1.9×) | 107 s (3.6×) | 62 s (6.1×) | 41 s (9.3×) |
| analytics-mts | 3.sh (vehicle hours on road) | 427 s | 232 s (1.8×) | 126 s (3.4×) | 74 s (5.8×) | 51 s (8.4×) |
| analytics-mts | 4.sh (hours monitored per day) | 167 s | 94 s (1.8×) | 54 s (3.1×) | 35 s (4.7×) | 28 s (6.0×) |
| oneliners | bi-grams.sh | 1007 s | 539 s (1.9×) | 286 s (3.5×) | 166 s (6.1×) | 118 s (8.6×) |
| oneliners | diff.sh | 478 s | 276 s (1.7×) | 167 s (2.9×) | 120 s (4.0×) | 98 s (4.9×) |
| oneliners | nfa-regex.sh | 391 s | 197 s (2.0×) | 99 s (3.9×) | 51 s (7.7×) | 26 s (14.9×) |
| oneliners | set-diff.sh | 1308 s | 717 s (1.8×) | 376 s (3.5×) | 220 s (6.0×) | 144 s (9.1×) |
| oneliners | shortest-scripts.sh | 110 s | 57 s (1.9×) | 29 s (3.8×) | 16 s (7.0×) | 9 s (12.4×) |
| oneliners | sort-sort.sh | 167 s | 101 s (1.7×) | 60 s (2.8×) | 40 s (4.2×) | 31 s (5.4×) |
| oneliners | sort.sh | 389 s | 207 s (1.9×) | 106 s (3.7×) | 59 s (6.6×) | 39 s (10.0×) |
| oneliners | spell.sh | 736 s | 386 s (1.9×) | 208 s (3.5×) | 115 s (6.4×) | 78 s (9.5×) |
| oneliners | top-n.sh | 622 s | 328 s (1.9×) | 169 s (3.7×) | 99 s (6.3×) | 63 s (9.9×) |
| oneliners | wf.sh | 2089 s | 1065 s (2.0×) | 545 s (3.8×) | 298 s (7.0×) | 196 s (10.7×) |
| poets | 1_1.sh (count_words) | 637 s | 443 s (1.4×) | 224 s (2.8×) | 123 s (5.2×) | 84 s (7.6×) |
| poets | 2_1.sh (merge_upper) | 547 s | 380 s (1.4×) | 195 s (2.8×) | 114 s (4.8×) | 79 s (6.9×) |
| poets | 2_2.sh (count_vowel_seq) | 140 s | 86 s (1.6×) | 52 s (2.7×) | 35 s (4.0×) | 27 s (5.2×) |
| poets | 3_1.sh (sort) | 665 s | 455 s (1.5×) | 232 s (2.9×) | 129 s (5.1×) | 89 s (7.4×) |
| poets | 3_2.sh (sort_words_by_folding) | 681 s | 467 s (1.5×) | 240 s (2.8×) | 136 s (5.0×) | 94 s (7.3×) |
| poets | 3_3.sh (sort_words_by_rhyming) | 699 s | 478 s (1.5×) | 246 s (2.8×) | 140 s (5.0×) | 100 s (7.0×) |
| poets | 4_3.sh (bigrams) | 915 s | 635 s (1.4×) | 346 s (2.6×) | 215 s (4.2×) | 173 s (5.3×) |
| poets | 4_3b.sh (count_trigrams) | 1049 s | 734 s (1.4×) | 430 s (2.4×) | 311 s (3.4×) | 275 s (3.8×) |
| poets | 6_1.sh (trigram_rec) | 5 s | 7 s (0.6×) | 7 s (0.7×) | 6 s (0.8×) | 8 s (0.6×) |
| poets | 6_1_1.sh (uppercase_by_token) | 45 s | 33 s (1.3×) | 22 s (2.1×) | 16 s (2.7×) | 14 s (3.3×) |
| poets | 6_1_2.sh (uppercase_by_type) | 635 s | 387 s (1.6×) | 188 s (3.4×) | 99 s (6.4×) | 64 s (10.0×) |
| poets | 6_2.sh (4letter_words) | 647 s | 399 s (1.6×) | 199 s (3.2×) | 108 s (6.0×) | 80 s (8.1×) |
| poets | 6_3.sh (words_no_vowels) | 235 s | 156 s (1.5×) | 83 s (2.8×) | 48 s (4.9×) | 32 s (7.4×) |
| poets | 6_4.sh (1syllable_words) | 542 s | 318 s (1.7×) | 164 s (3.3×) | 91 s (6.0×) | 57 s (9.5×) |
| poets | 6_5.sh (2syllable_words) | 443 s | 282 s (1.6×) | 143 s (3.1×) | 78 s (5.7×) | 48 s (9.2×) |
| poets | 6_7.sh (verses_2om_3om_2instances) | 7 s | 11 s (0.7×) | 10 s (0.7×) | 10 s (0.7×) | 11 s (0.6×) |
| poets | 7_2.sh (count_consonant_seq) | 678 s | 370 s (1.8×) | 198 s (3.4×) | 119 s (5.7×) | 80 s (8.5×) |
| poets | 8.2_1.sh (vowel_sequences_gr_1K) | 573 s | 348 s (1.6×) | 186 s (3.1×) | 110 s (5.2×) | 73 s (7.9×) |
| poets | 8.2_2.sh (bigrams_appear_twice) | 921 s | 642 s (1.4×) | 351 s (2.6×) | 222 s (4.2×) | 177 s (5.2×) |
| poets | 8.3_2.sh (find_anagrams) | 724 s | 440 s (1.6×) | 227 s (3.2×) | 133 s (5.4×) | 102 s (7.1×) |
| poets | 8.3_3.sh (compare_exodus_genesis) | 656 s | 401 s (1.6×) | 197 s (3.3×) | 108 s (6.1×) | 74 s (8.8×) |
| poets | 8_1.sh (sort_words_by_n_syllables) | 653 s | 403 s (1.6×) | 196 s (3.3×) | 104 s (6.3×) | 69 s (9.5×) |
| unix50 | 1.sh (1.0: extract last name) | 12 s | 8 s (1.5×) | 5 s (2.4×) | 4 s (3.2×) | 3 s (3.6×) |
| unix50 | 10.sh (4.4: histogram by piece) | 48 s | 32 s (1.5×) | 19 s (2.5×) | 14 s (3.4×) | 13 s (3.7×) |
| unix50 | 11.sh (4.5: histogram by piece and pawn) | 42 s | 28 s (1.5×) | 18 s (2.4×) | 13 s (3.1×) | 12 s (3.4×) |
| unix50 | 12.sh (4.6: piece used most) | 149 s | 91 s (1.6×) | 53 s (2.8×) | 35 s (4.3×) | 27 s (5.5×) |
| unix50 | 13.sh (5.1: extract hellow world) | 12 s | 10 s (1.2×) | 7 s (1.7×) | 6 s (2.0×) | 7 s (1.7×) |
| unix50 | 14.sh (6.1: order bodies) | 185 s | 106 s (1.7×) | 61 s (3.0×) | 40 s (4.6×) | 31 s (6.0×) |
| unix50 | 15.sh (7.1: number of versions) | 8 s | 7 s (1.2×) | 5 s (1.6×) | 5 s (1.6×) | 6 s (1.4×) |
| unix50 | 16.sh (7.2: most frequent machine) | 93 s | 53 s (1.7×) | 30 s (3.1×) | 19 s (4.8×) | 15 s (6.4×) |
| unix50 | 17.sh (7.3: decades unix released) | 43 s | 26 s (1.6×) | 17 s (2.6×) | 12 s (3.7×) | 10 s (4.4×) |
| unix50 | 18.sh (8.1: count unix birth-year) | 3 s | 5 s (0.6×) | 4 s (0.7×) | 5 s (0.6×) | 6 s (0.5×) |
| unix50 | 19.sh (8.2: location office) | 2 s | 3 s (0.9×) | 2 s (1.0×) | 2 s (1.0×) | 2 s (1.0×) |
| unix50 | 2.sh (1.1: extract names and sort) | 171 s | 98 s (1.7×) | 51 s (3.3×) | 31 s (5.6×) | 20 s (8.4×) |
| unix50 | 20.sh (8.3: four most involved) | 5 s | 6 s (0.8×) | 5 s (1.0×) | 5 s (1.0×) | 6 s (0.8×) |
| unix50 | 21.sh (8.4: longest words w/o hyphens) | 733 s | 384 s (1.9×) | 192 s (3.8×) | 104 s (7.0×) | 64 s (11.4×) |
| unix50 | 23.sh (9.1: extract word PORT) | 202 s | 109 s (1.9×) | 61 s (3.3×) | 35 s (5.7×) | 23 s (8.8×) |
| unix50 | 24.sh (9.2: extract word BELL) | 5 s | 4 s (1.3×) | 3 s (1.8×) | 2 s (2.0×) | 2 s (2.1×) |
| unix50 | 25.sh (9.3: animal decorate) | 6 s | 4 s (1.3×) | 3 s (1.7×) | 3 s (2.0×) | 3 s (2.1×) |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| unix50 | 26.sh (9.4: four corners) | 32 s | 26 s (1.2×) | 21 s (1.5×) | 18 s (1.7×) | 18 s (1.7×) |
| unix50 | 28.sh (9.6: follow directions) | 188 s | 119 s (1.6×) | 82 s (2.3×) | 64 s (2.9×) | 54 s (3.5×) |
| unix50 | 29.sh (9.7: four corners) | 19 s | 20 s (0.9×) | 18 s (1.1×) | 18 s (1.0×) | 18 s (1.0×) |
| unix50 | 3.sh (1.2: extract names and sort) | 0 s | 0 s (1.0×) | 0 s (0.7×) | 0 s (0.7×) | 0 s (0.7×) |
| unix50 | 30.sh (9.8: TELE-communications) | 154 s | 119 s (1.3×) | 85 s (1.8×) | 72 s (2.1×) | 66 s (2.3×) |
| unix50 | 31.sh (9.9) | 149 s | 111 s (1.3×) | 89 s (1.7×) | 78 s (1.9×) | 73 s (2.0×) |
| unix50 | 32.sh (10.1: count recipients) | 6 s | 6 s (0.9×) | 6 s (1.0×) | 6 s (1.0×) | 7 s (0.9×) |
| unix50 | 33.sh (10.2: list recipients) | 6 s | 6 s (1.0×) | 5 s (1.1×) | 6 s (1.1×) | 6 s (1.0×) |
| unix50 | 34.sh (10.3: extract username) | 2 s | 3 s (0.9×) | 2 s (1.0×) | 3 s (0.9×) | 3 s (0.7×) |
| unix50 | 35.sh (11.1: year received medal) | 1 s | 2 s (0.5×) | 2 s (0.6×) | 2 s (0.5×) | 2 s (0.5×) |
| unix50 | 36.sh (11.2: most repeated first name) | 19 s | 12 s (1.6×) | 8 s (2.4×) | 6 s (3.1×) | 7 s (2.8×) |
| unix50 | 4.sh (1.3: sort top first names) | 154 s | 89 s (1.7×) | 49 s (3.2×) | 30 s (5.1×) | 21 s (7.4×) |
| unix50 | 5.sh (2.1: all Unix utilities) | 8 s | 5 s (1.4×) | 4 s (2.1×) | 3 s (2.7×) | 3 s (2.5×) |
| unix50 | 6.sh (3.1: first letter of last names) | 14 s | 10 s (1.4×) | 7 s (2.1×) | 5 s (2.6×) | 5 s (2.7×) |
| unix50 | 7.sh (4.1: number of rounds) | 18 s | 14 s (1.3×) | 10 s (1.8×) | 9 s (2.1×) | 9 s (2.0×) |
| unix50 | 8.sh (4.2: pieces captured) | 12 s | 11 s (1.2×) | 8 s (1.6×) | 7 s (1.7×) | 8 s (1.6×) |
| unix50 | 9.sh (4.3: pieces captured with pawn) | 28 s | 19 s (1.5×) | 12 s (2.2×) | 10 s (2.9×) | 9 s (2.9×) |
| **Max** | | 2089 s | 1065 s (2.0×) | 545 s (3.9×) | 311 s (7.7×) | 275 s (14.9×) |
| **Min** | | 0 s | 0 s (0.5×) | 0 s (0.6×) | 0 s (0.5×) | 0 s (0.5×) |
| **Mean** | | 332 s | 200 s (1.5×) | 107 s (2.5×) | 65 s (4.0×) | 48 s (5.5×) |
| **Median** | | 167 s | 104 s (1.5×) | 60 s (2.8×) | 38 s (4.2×) | 28 s (5.3×) |

**Table 6.** Performance results for all benchmark scripts, where new pipelines are optimized by eliminating intermediate combiners

| Benchmark | Script Name | $u_1$ | $T_1$ | $T_2$ | $T_4$ | $T_8$ | $T_{16}$ |
|---|---|---|---|---|---|---|---|
| analytics-mts | 1.sh (vehicles per day) | 376 s | 330 s (1.1×) | 170 s (2.2×) | 89 s (4.2×) | 48 s (7.9×) | 29 s (13.1×) |
| analytics-mts | 2.sh (vehicle days on road) | 379 s | 331 s (1.1×) | 169 s (2.2×) | 88 s (4.3×) | 47 s (8.0×) | 28 s (13.5×) |
| analytics-mts | 3.sh (vehicle hours on road) | 427 s | 411 s (1.0×) | 213 s (2.0×) | 112 s (3.8×) | 61 s (6.9×) | 38 s (11.3×) |
| analytics-mts | 4.sh (hours monitored per day) | 167 s | 97 s (1.7×) | 53 s (3.2×) | 29 s (5.8×) | 18 s (9.4×) | 13 s (12.8×) |
| oneliners | bi-grams.sh | 1007 s | 1015 s (1.0×) | 535 s (1.9×) | 283 s (3.6×) | 168 s (6.0×) | 115 s (8.7×) |
| oneliners | diff.sh | 478 s | 332 s (1.4×) | 226 s (2.1×) | 137 s (3.5×) | 100 s (4.8×) | 83 s (5.8×) |
| oneliners | nfa-regex.sh | 391 s | 388 s (1.0×) | 196 s (2.0×) | 99 s (3.9×) | 51 s (7.7×) | 27 s (14.7×) |
| oneliners | set-diff.sh | 1308 s | 816 s (1.6×) | 495 s (2.6×) | 279 s (4.7×) | 175 s (7.5×) | 128 s (10.2×) |
| oneliners | shortest-scripts.sh | 110 s | 82 s (1.3×) | 42 s (2.6×) | 22 s (5.0×) | 12 s (9.5×) | 7 s (16.2×) |
| oneliners | sort-sort.sh | 167 s | 137 s (1.2×) | 85 s (2.0×) | 53 s (3.2×) | 35 s (4.7×) | 28 s (6.0×) |
| oneliners | sort.sh | 389 s | 391 s (1.0×) | 207 s (1.9×) | 106 s (3.7×) | 59 s (6.6×) | 38 s (10.3×) |
| oneliners | spell.sh | 736 s | 484 s (1.5×) | 282 s (2.6×) | 154 s (4.8×) | 90 s (8.1×) | 61 s (12.1×) |
| oneliners | top-n.sh | 622 s | 388 s (1.6×) | 228 s (2.7×) | 127 s (4.9×) | 75 s (8.3×) | 50 s (12.4×) |
| oneliners | wf.sh | 2089 s | 1196 s (1.7×) | 667 s (3.1×) | 368 s (5.7×) | 223 s (9.4×) | 145 s (14.4×) |
| poets | 1_1.sh (count_words) | 637 s | 637 s (1.0×) | 440 s (1.4×) | 223 s (2.9×) | 123 s (5.2×) | 83 s (7.6×) |
| poets | 2_1.sh (merge_upper) | 547 s | 543 s (1.0×) | 376 s (1.5×) | 193 s (2.8×) | 109 s (5.0×) | 78 s (7.0×) |
| poets | 2_2.sh (count_vowel_seq) | 140 s | 142 s (1.0×) | 83 s (1.7×) | 50 s (2.8×) | 40 s (3.5×) | 24 s (5.8×) |
| poets | 3_1.sh (sort) | 665 s | 670 s (1.0×) | 460 s (1.4×) | 232 s (2.9×) | 132 s (5.0×) | 88 s (7.6×) |
| poets | 3_2.sh (sort_words_by_folding) | 681 s | 685 s (1.0×) | 472 s (1.4×) | 237 s (2.9×) | 134 s (5.1×) | 94 s (7.2×) |
| poets | 3_3.sh (sort_words_by_rhyming) | 699 s | 704 s (1.0×) | 478 s (1.5×) | 250 s (2.8×) | 140 s (5.0×) | 100 s (7.0×) |
| poets | 4_3.sh (bigrams) | 915 s | 909 s (1.0×) | 640 s (1.4×) | 343 s (2.7×) | 216 s (4.2×) | 173 s (5.3×) |
| poets | 4_3b.sh (count_trigrams) | 1049 s | 1056 s (1.0×) | 733 s (1.4×) | 430 s (2.4×) | 311 s (3.4×) | 279 s (3.8×) |
| poets | 6_1.sh (trigram_rec) | 5 s | 2 s (1.9×) | 2 s (2.8×) | 2 s (2.4×) | 1 s (3.6×) | 2 s (2.4×) |
| poets | 6_1_1.sh (uppercase_by_token) | 45 s | 40 s (1.1×) | 30 s (1.5×) | 21 s (2.1×) | 16 s (2.8×) | 14 s (3.2×) |
| poets | 6_1_2.sh (uppercase_by_type) | 635 s | 158 s (4.0×) | 101 s (6.3×) | 53 s (12.1×) | 31 s (20.5×) | 24 s (26.9×) |
| poets | 6_2.sh (4letter_words) | 647 s | 171 s (3.8×) | 114 s (5.7×) | 63 s (10.2×) | 41 s (15.6×) | 34 s (18.8×) |
| poets | 6_3.sh (words_no_vowels) | 235 s | 223 s (1.1×) | 148 s (1.6×) | 82 s (2.9×) | 46 s (5.1×) | 31 s (7.7×) |
| poets | 6_4.sh (1syllable_words) | 542 s | 203 s (2.7×) | 139 s (3.9×) | 76 s (7.1×) | 45 s (11.9×) | 31 s (17.4×) |
| poets | 6_5.sh (2syllable_words) | 443 s | 316 s (1.4×) | 205 s (2.2×) | 109 s (4.1×) | 61 s (7.3×) | 40 s (11.0×) |
| poets | 6_7.sh (verses_2om_3om_2instances) | 7 s | 4 s (1.8×) | 3 s (2.8×) | 3 s (2.7×) | 3 s (2.4×) | 5 s (1.5×) |
| poets | 7_2.sh (count_consonant_seq) | 678 s | 250 s (2.7×) | 152 s (4.5×) | 88 s (7.7×) | 59 s (11.5×) | 48 s (14.2×) |
| poets | 8.2_1.sh (vowel_sequences_gr_1K) | 573 s | 155 s (3.7×) | 101 s (5.7×) | 66 s (8.7×) | 50 s (11.4×) | 42 s (13.7×) |
| poets | 8.2_2.sh (bigrams_appear_twice) | 921 s | 247 s (3.7×) | 182 s (5.0×) | 120 s (7.7×) | 96 s (9.6×) | 91 s (10.2×) |
| poets | 8.3_2.sh (find_anagrams) | 724 s | 208 s (3.5×) | 131 s (5.5×) | 76 s (9.5×) | 55 s (13.3×) | 50 s (14.5×) |
| poets | 8.3_3.sh (compare_exodus_genesis) | 656 s | 162 s (4.1×) | 106 s (6.2×) | 59 s (11.1×) | 39 s (16.7×) | 34 s (19.3×) |

| | | | | | | |
|---|---|---|---|---|---|---|
| poets | 8_1.sh (sort_words_by_n syllables) | 653 s | 169 s (3.9×) | 108 s (6.0×) | 58 s (11.3×) | 35 s (18.4×) | 26 s (24.6×) |
| unix50 | 1.sh (1.0: extract last name) | 12 s | 12 s (1.0×) | 8 s (1.5×) | 5 s (2.4×) | 4 s (3.2×) | 3 s (3.5×) |
| unix50 | 10.sh (4.4: histogram by piece) | 48 s | 27 s (1.7×) | 17 s (2.9×) | 10 s (4.9×) | 7 s (6.7×) | 6 s (7.8×) |
| unix50 | 11.sh (4.5: histogram by piece and pawn) | 42 s | 24 s (1.7×) | 16 s (2.6×) | 10 s (4.3×) | 7 s (5.9×) | 6 s (6.8×) |
| unix50 | 12.sh (4.6: piece used most) | 149 s | 112 s (1.3×) | 69 s (2.2×) | 38 s (4.0×) | 24 s (6.2×) | 18 s (8.2×) |
| unix50 | 13.sh (5.1: extract hellow world) | 12 s | 5 s (2.6×) | 4 s (3.1×) | 3 s (4.1×) | 3 s (4.7×) | 2 s (5.1×) |
| unix50 | 14.sh (6.1: order bodies) | 185 s | 154 s (1.2×) | 92 s (2.0×) | 51 s (3.6×) | 33 s (5.6×) | 25 s (7.5×) |
| unix50 | 15.sh (7.1: number of versions) | 8 s | 6 s (1.4×) | 4 s (2.0×) | 3 s (2.9×) | 3 s (3.0×) | 3 s (2.5×) |
| unix50 | 16.sh (7.2: most frequent machine) | 93 s | 85 s (1.1×) | 48 s (2.0×) | 27 s (3.4×) | 17 s (5.4×) | 13 s (7.4×) |
| unix50 | 17.sh (7.3: decades unix released) | 43 s | 41 s (1.1×) | 24 s (1.8×) | 15 s (3.0×) | 10 s (4.2×) | 8 s (5.1×) |
| unix50 | 18.sh (8.1: count unix birth-year) | 3 s | 2 s (1.5×) | 2 s (1.4×) | 2 s (1.6×) | 2 s (1.4×) | 3 s (1.1×) |
| unix50 | 19.sh (8.2: location office) | 2 s | 2 s (1.1×) | 2 s (1.0×) | 2 s (1.2×) | 2 s (1.4×) | 2 s (1.0×) |
| unix50 | 2.sh (1.1: extract names and sort) | 171 s | 131 s (1.3×) | 74 s (2.3×) | 41 s (4.1×) | 25 s (6.8×) | 17 s (9.9×) |
| unix50 | 20.sh (8.3: four most involved) | 5 s | 0 s (NaN) | 1 s (3.9×) | 1 s (3.6×) | 2 s (2.5×) | 3 s (1.7×) |
| unix50 | 21.sh (8.4: longest words w/o hyphens) | 733 s | 440 s (1.7×) | 257 s (2.8×) | 141 s (5.2×) | 78 s (9.4×) | 49 s (14.9×) |
| unix50 | 23.sh (9.1: extract word PORT) | 202 s | 116 s (1.7×) | 59 s (3.4×) | 31 s (6.5×) | 17 s (11.7×) | 10 s (19.8×) |
| unix50 | 24.sh (9.2: extract word BELL) | 5 s | 5 s (1.1×) | 4 s (1.4×) | 2 s (2.1×) | 2 s (2.3×) | 2 s (2.4×) |
| unix50 | 25.sh (9.3: animal decorate) | 6 s | 5 s (1.1×) | 4 s (1.5×) | 3 s (2.0×) | 2 s (2.6×) | 2 s (2.3×) |
| unix50 | 26.sh (9.4: four corners) | 32 s | 28 s (1.1×) | 22 s (1.5×) | 17 s (1.9×) | 16 s (2.0×) | 15 s (2.1×) |
| unix50 | 28.sh (9.6: follow directions) | 188 s | 185 s (1.0×) | 114 s (1.7×) | 78 s (2.4×) | 59 s (3.2×) | 49 s (3.8×) |
| unix50 | 29.sh (9.7: four corners) | 19 s | 16 s (1.2×) | 16 s (1.2×) | 15 s (1.2×) | 15 s (1.3×) | 15 s (1.3×) |
| unix50 | 3.sh (1.2: extract names and sort) | 0 s | 0 s (0.9×) | 0 s (1.0×) | 0 s (0.7×) | 0 s (0.7×) | 0 s (0.7×) |
| unix50 | 30.sh (9.8: TELE-communications) | 154 s | 152 s (1.0×) | 105 s (1.5×) | 80 s (1.9×) | 71 s (2.2×) | 62 s (2.5×) |
| unix50 | 31.sh (9.9) | 149 s | 145 s (1.0×) | 106 s (1.4×) | 85 s (1.8×) | 74 s (2.0×) | 68 s (2.2×) |
| unix50 | 32.sh (10.1: count recipients) | 6 s | 5 s (1.1×) | 6 s (1.0×) | 5 s (1.1×) | 6 s (1.0×) | 6 s (0.9×) |
| unix50 | 33.sh (10.2: list recipients) | 6 s | 6 s (1.0×) | 6 s (1.0×) | 5 s (1.1×) | 5 s (1.1×) | 5 s (1.1×) |
| unix50 | 34.sh (10.3: extract username) | 2 s | 0 s (19.7×) | 1 s (1.9×) | 1 s (1.8×) | 2 s (1.3×) | 3 s (0.9×) |
| unix50 | 35.sh (11.1: year received medal) | 1 s | 1 s (1.1×) | 1 s (0.6×) | 1 s (0.6×) | 2 s (0.6×) | 2 s (0.6×) |
| unix50 | 36.sh (11.2: most repeated first name) | 19 s | 15 s (1.3×) | 10 s (1.9×) | 7 s (2.9×) | 5 s (3.6×) | 6 s (3.5×) |
| unix50 | 4.sh (1.3: sort top first names) | 154 s | 131 s (1.2×) | 79 s (2.0×) | 43 s (3.6×) | 26 s (5.9×) | 19 s (8.3×) |
| unix50 | 5.sh (2.1: all Unix utilities) | 8 s | 7 s (1.1×) | 5 s (1.6×) | 3 s (2.5×) | 2 s (3.2×) | 2 s (3.1×) |
| unix50 | 6.sh (3.1: first letter of last names) | 14 s | 10 s (1.4×) | 7 s (2.2×) | 4 s (3.4×) | 3 s (4.1×) | 3 s (5.2×) |
| unix50 | 7.sh (4.1: number of rounds) | 18 s | 14 s (1.3×) | 9 s (2.0×) | 5 s (3.5×) | 4 s (4.7×) | 4 s (4.7×) |
| unix50 | 8.sh (4.2: pieces captured) | 12 s | 6 s (2.1×) | 4 s (3.0×) | 3 s (4.2×) | 3 s (5.0×) | 3 s (4.1×) |
| unix50 | 9.sh (4.3: pieces captured with pawn) | 28 s | 14 s (2.0×) | 8 s (3.3×) | 5 s (5.6×) | 4 s (7.5×) | 4 s (7.3×) |
| **Max** | | 2089 s | 1196 s (19.7×) | 733 s (6.3×) | 430 s (12.1×) | 311 s (20.5×) | 279 s (26.9×) |
| **Min** | | 0 s | 0 s (0.9×) | 0 s (0.6×) | 0 s (0.6×) | 0 s (0.6×) | 0 s (0.6×) |
| **Mean** | | 332 s | 228 s (1.9×) | 142 s (2.4×) | 79 s (4.0×) | 50 s (6.1×) | 37 s (8.0×) |
| **Median** | | 167 s | 140 s (1.2×) | 84 s (2.0×) | 50 s (3.5×) | 32 s (5.1×) | 24 s (7.1×) |

# D  Appendix: Combiner Synthesis Results

Table 8 summarizes all plausible combiners identified by KumQuat during the combiner synthesis for the benchmarks. The first column presents the number of times the combiner appears as plausible across all benchmark scripts. The second column presents the combiner in the KumQuat combiner DSL, with variables $a, b$ denoting input arguments and $*$ denoting any flags. For each command in the benchmarks, the synthesized plausible combiners are all equivalent when operating on the command's outputs.

Table 9 presents all of the benchmark commands for which KumQuat did not synthesize a combiner.

Table 10 presents comprehensive synthesis results from our benchmark set of commands. The first column presents the name of the command $f$. The second presents the command itself. The third presents the size of the candidate combiners with a finite scope of eight. Inside parentheses is the breakdown of the candidates into the three combiner classes RecOp, StructOp, RunOp$_f$. The fourth column presents the wall-clock synthesis time in seconds. The fifth column presents the final set of synthesized plausible combiners expressed in the combiner DSL, where variables $a, b$ denote the first and the second input streams, respectively. The last column presents the number of these plausible combiners.

**Table 10.** Synthesis results for unique command/flag combinations

| Bench | Script | Idx | Command | Search Space | Time | Synthesized Plausible | #P |
|---|---|---|---|---|---|---|---|

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| oneliners | spell | 8 | `IN=${IN:-../benchmarks/pipelines/one liners/input/1G.txt} dict=${dict:-../in/dict.sorted} LC_COLLATE=C comm -23 - $dict` | 26404 (= 12440 + 13960 + 4) | 331 s | $e_1$ = (**concat** $a$ $b$), $e_2$ = (**rerun** $a$ $b$). | 2 |
| poets | vowel_-sequencies_-gr_1K | 7 | `awk "\$1 >= 1000"` | 26404 (= 12440 + 13960 + 4) | 176 s | $e_1$ = (**concat** $a$ $b$), $e_2$ = (**rerun** $a$ $b$). | 2 |
| poets | find_-anagrams | 8 | `awk "\$1 >= 2 {print \$2}"` | 2700 (= 968 + 1728 + 4) | 40 s | $e_1$ = (**concat** $a$ $b$). | 1 |
| unix50 | 8.4: longest words w/o hyphens | 3 | `awk "length >= 16"` | 26404 (= 12440 + 13960 + 4) | 172 s | $e_1$ = (**concat** $a$ $b$), $e_2$ = (**rerun** $a$ $b$). | 2 |
| unix50 | 8.2: location office | 4 | `awk "{\$1=\$1};1"` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1$ = (**concat** $a$ $b$), $e_2$ = (**rerun** $a$ $b$). | 2 |
| unix50 | 6.1: order bodies | 1 | `awk "{print \$2, \$0}"` | 110444 (= 59048 + 51392 + 4) | 124 s | $e_1$ = (**concat** $a$ $b$). | 1 |
| unix50 | 8.2: location office | 2 | `awk 'length <= 45'` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1$ = (**concat** $a$ $b$), $e_2$ = (**rerun** $a$ $b$). | 2 |
| poets | sort_-words_-by_n_-syllables | 6 | `awk '{print NF}'` | 2700 (= 968 + 1728 + 4) | 39 s | $e_1$ = (**concat** $a$ $b$). | 1 |
| analytics-mts | vehicles per day | 7 | `awk -v OFS="\t" "{print \$2,\$1}"` | 110444 (= 59048 + 51392 + 4) | 125 s | $e_1$ = (**concat** $a$ $b$). | 1 |
| analytics-mts | vehicles per day | 0 | `cat` | 26404 (= 12440 + 13960 + 4) | 59 s | $e_1$ = (**concat** $a$ $b$), $e_2$ = (**rerun** $a$ $b$). | 2 |
| oneliners | spell | 2 | `col -bx` | 26404 (= 12440 + 13960 + 4) | 59 s | $e_1$ = (**concat** $a$ $b$), $e_2$ = (**rerun** $a$ $b$). | 2 |
| unix50 | 4.4: histogram by piece | 6 | `cut -c 1-1` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1$ = (**concat** $a$ $b$), $e_2$ = (**rerun** $a$ $b$). | 2 |
| unix50 | 5.1: extract hellow world | 3 | `cut -c 1-12` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1$ = (**concat** $a$ $b$), $e_2$ = (**rerun** $a$ $b$). | 2 |
| unix50 | 9.3: animal decorate | 1 | `cut -c 1-2` | 26404 (= 12440 + 13960 + 4) | 59 s | $e_1$ = (**concat** $a$ $b$), $e_2$ = (**rerun** $a$ $b$). | 2 |
| unix50 | 9.1: extract word PORT | 6 | `cut -c 1-4` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1$ = (**concat** $a$ $b$), $e_2$ = (**rerun** $a$ $b$). | 2 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| unix50 | 7.3: decades unix released | 3 | `cut -c 3-3` | 26404 (= 12440 + 13960 + 4) | 59 s | $e_1 = (\textbf{concat } a\ b)$. | 1 |
| unix50 | 5.1: extract hellow world | 2 | `cut -d "\"" -f 2` | 26404 (= 12440 + 13960 + 4) | 104 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| oneliners | set-diff | 2 | `cut -d ' ' -f 1` | 2700 (= 968 + 1728 + 4) | 71 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| unix50 | 1.0: extract last name | 1 | `cut -d ' ' -f 2` | 2700 (= 968 + 1728 + 4) | 72 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| unix50 | 2.1: all Unix utilities | 1 | `cut -d ' ' -f 4` | 2700 (= 968 + 1728 + 4) | 71 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| unix50 | 8.3: four most involved | 2 | `cut -d '(' -f 2` | 26404 (= 12440 + 13960 + 4) | 102 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| unix50 | 8.3: four most involved | 3 | `cut -d ')' -f 1` | 26404 (= 12440 + 13960 + 4) | 102 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| analytics-mts | vehicles per day | 4 | `cut -d ',' -f 1` | 26404 (= 12440 + 13960 + 4) | 102 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| analytics-mts | hours monitored per day | 2 | `cut -d ',' -f 1,2` | 110444 (= 59048 + 51392 + 4) | 126 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| analytics-mts | vehicle hours on road | 2 | `cut -d ',' -f 1,2,4` | 110444 (= 59048 + 51392 + 4) | 126 s | $e_1 = (\textbf{concat } a\ b)$. | 1 |
| analytics-mts | vehicles per day | 2 | `cut -d ',' -f 1,3` | 110444 (= 59048 + 51392 + 4) | 125 s | $e_1 = (\textbf{concat } a\ b)$. | 1 |
| analytics-mts | vehicle days on road | 4 | `cut -d ',' -f 2` | 26404 (= 12440 + 13960 + 4) | 103 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| analytics-mts | vehicle hours on road | 4 | `cut -d ',' -f 3` | 26404 (= 12440 + 13960 + 4) | 101 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| analytics-mts | vehicle days on road | 2 | `cut -d ',' -f 3,1` | 110444 (= 59048 + 51392 + 4) | 124 s | $e_1 = (\textbf{concat } a\ b)$. | 1 |
| unix50 | 4.4: histogram by piece | 4 | `cut -d '.' -f 2` | 26404 (= 12440 + 13960 + 4) | 103 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |

| oneliners | shortest-scripts | 3 | `cut -d: -f1` | 26404 (= 12440 + 13960 + 4) | 61 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |
| --- | --- | --- | --- | --- | --- | --- | --- |
| unix50 | 7.1: number of versions | 1 | `cut -f 1` | 26404 (= 12440 + 13960 + 4) | 102 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |
| unix50 | 7.2: most frequent machine | 1 | `cut -f 2` | 26404 (= 12440 + 13960 + 4) | 103 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |
| unix50 | 7.3: decades unix released | 1 | `cut -f 4` | 26404 (= 12440 + 13960 + 4) | 102 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |
| unix50 | 10.3: extract username | 4 | `fmt -w1` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |
| unix50 | 9.4: four corners | 2 | `grep "\""` | 26404 (= 12440 + 13960 + 4) | 61 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |
| oneliners | shortest-scripts | 2 | `grep "shell script"` | 26404 (= 12440 + 13960 + 4) | 62 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |
| unix50 | 8.3: four most involved | 1 | `grep '('` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |
| unix50 | 7.1: number of versions | 2 | `grep 'AT&T'` | 26404 (= 12440 + 13960 + 4) | 61 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |
| poets | trigram_-rec | 10 | `grep 'And he said'` | 26404 (= 12440 + 13960 + 4) | 61 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |
| unix50 | 8.2: location office | 1 | `grep 'Bell'` | 26404 (= 12440 + 13960 + 4) | 61 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |
| unix50 | 11.1: year received medal | 1 | `grep 'UNIX'` | 26404 (= 12440 + 13960 + 4) | 61 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |
| unix50 | 9.1: extract word PORT | 2 | `grep '[A-Z]'` | 26404 (= 12440 + 13960 + 4) | 59 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |
| unix50 | 4.4: histogram by piece | 5 | `grep '[KQRBN]'` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat}\ a\ b),$ $e_2 = (\textbf{rerun}\ a\ b).$ | 2 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| oneliners | nfa-regex | 2 | `grep '\(.\).*\1\(.\).*\2\(.\).*\3\(.\).*\4'` | 110444 (= 59048 + 51392 + 4) | 129 s | $e_1 = (\textbf{concat } a\ b),$ $e_2 = (\textbf{rerun } a\ b).$ | 2 |
| unix50 | 4.4: histogram by piece | 3 | `grep '\.'` | 26404 (= 12440 + 13960 + 4) | 165 s | $e_1 = (\textbf{concat } a\ b),$ $e_2 = (\textbf{rerun } a\ b).$ | 2 |
| poets | verses_-2om_-3om_-2instances | 11 | `grep 'light.\*light'` | 26404 (= 12440 + 13960 + 4) | 172 s | $e_1 = (\textbf{concat } a\ b),$ $e_2 = (\textbf{rerun } a\ b).$ | 2 |
| unix50 | 5.1: extract hellow world | 1 | `grep 'print'` | 26404 (= 12440 + 13960 + 4) | 61 s | $e_1 = (\textbf{concat } a\ b),$ $e_2 = (\textbf{rerun } a\ b).$ | 2 |
| poets | trigram_-rec | 3 | `grep 'the land of'` | 26404 (= 12440 + 13960 + 4) | 61 s | $e_1 = (\textbf{concat } a\ b),$ $e_2 = (\textbf{rerun } a\ b).$ | 2 |
| unix50 | 4.4: histogram by piece | 2 | `grep 'x'` | 26404 (= 12440 + 13960 + 4) | 179 s | $e_1 = (\textbf{concat } a\ b),$ $e_2 = (\textbf{rerun } a\ b).$ | 2 |
| poets | 4letter_-words | 10 | `grep -c '^....$'` | 2700 (= 968 + 1728 + 4) | 38 s | $e_1 = ((\textbf{back } `\backslash n` \textbf{ add})\ a\ b),$ $e_2 = ((\textbf{back } `\backslash n` \textbf{ add})\ b\ a).$ | 2 |
| poets | uppercase_-by_-token | 4 | `grep -c '^[A-Z]'` | 2700 (= 968 + 1728 + 4) | 39 s | $e_1 = ((\textbf{back } `\backslash n` \textbf{ add})\ a\ b),$ $e_2 = ((\textbf{back } `\backslash n` \textbf{ add})\ b\ a).$ | 2 |
| poets | verses_-2om_-3om_-2instances | 3 | `grep -c 'light.\*light'` | 2700 (= 968 + 1728 + 4) | 39 s | $e_1 = ((\textbf{back } `\backslash n` \textbf{ add})\ a\ b),$ $e_2 = ((\textbf{back } `\backslash n` \textbf{ add})\ b\ a).$ | 2 |
| poets | verses_-2om_-3om_-2instances | 7 | `grep -c 'light.\*light.\*light'` | 2700 (= 968 + 1728 + 4) | 39 s | $e_1 = ((\textbf{back } `\backslash n` \textbf{ add})\ a\ b),$ $e_2 = ((\textbf{back } `\backslash n` \textbf{ add})\ b\ a).$ | 2 |
| poets | 1syllable_-words | 4 | `grep -i '^[^aeiou]*[aeiou][^aeiou]*$'` | 110444 (= 59048 + 51392 + 4) | 177 s | $e_1 = (\textbf{concat } a\ b),$ $e_2 = (\textbf{merge } a\ b),$ $e_3 = (\textbf{merge } b\ a),$ $e_4 = (\textbf{rerun } a\ b).$ | 4 |
| poets | 2syllable_-words | 4 | `grep -i '^[^aeiou]*[aeiou][^aeiou]*[aeiou][^aeiou]$'` | 110444 (= 59048 + 51392 + 4) | 315 s | $e_1 = (\textbf{concat } a\ b),$ $e_2 = (\textbf{rerun } a\ b).$ | 2 |
| unix50 | 4.3: pieces captured with pawn | 5 | `grep -v '[KQRBN]'` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat } a\ b),$ $e_2 = (\textbf{rerun } a\ b).$ | 2 |
| oneliners | shortest-scripts | 5 | `grep -v '^0$'` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat } a\ b),$ $e_2 = (\textbf{rerun } a\ b).$ | 2 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| poets | verses_-2om_-3om_-2instances | 12 | `grep -vc 'light.\*light.\*light'` | 2700 $(= 968 + 1728 + 4)$ | 39 s | $e_1 = ((\textbf{back } `\backslash n` \textbf{ add}) \, a \, b),$ <br> $e_2 = ((\textbf{back } `\backslash n` \textbf{ add}) \, b \, a).$ | 2 |
| poets | words_-no_-vowels | 4 | `grep -vi '[aeiou]'` | 26404 $(= 12440 + 13960 + 4)$ | 59 s | $e_1 = (\textbf{concat } a \, b),$ <br> $e_2 = (\textbf{rerun } a \, b).$ | 2 |
| unix50 | 8.1: count unix birth-year | 2 | `grep 1969` | 26404 $(= 12440 + 13960 + 4)$ | 60 s | $e_1 = (\textbf{concat } a \, b),$ <br> $e_2 = (\textbf{rerun } a \, b).$ | 2 |
| poets | compare_-exodus_-genesis | 9 | `head` | 26404 $(= 12440 + 13960 + 4)$ | 163 s | $e_1 = (\textbf{rerun } a \, b).$ | 1 |
| oneliners | shortest-scripts | 7 | `head -15` | 26404 $(= 12440 + 13960 + 4)$ | 185 s | $e_1 = (\textbf{rerun } a \, b).$ | 1 |
| unix50 | 7.2: most fre-quent ma-chine | 5 | `head -n 1` | 26404 $(= 12440 + 13960 + 4)$ | 102 s | $e_1 = (\textbf{first } a \, b),$ <br> $e_2 = (\textbf{second } b \, a),$ <br> $e_3 = ((\textbf{back } `\backslash n` \textbf{ first}) \, a \, b),$ <br> $e_4 = ((\textbf{fuse } `\backslash n` \textbf{ first}) \, a \, b),$ <br> $e_5 = ((\textbf{back } `\backslash n` \textbf{ second}) \, b \, a),$ <br> $e_6 = ((\textbf{fuse } `\backslash n` \textbf{ second}) \, b \, a),$ <br> $e_7 = (\textbf{rerun } a \, b).$ | 7 |
| unix50 | 1.2: extract names and sort | 1 | `head -n 2` | 26404 $(= 12440 + 13960 + 4)$ | 101 s | $e_1 = (\textbf{rerun } a \, b).$ | 1 |
| unix50 | 4.6: piece used most | 8 | `head -n 3` | 26404 $(= 12440 + 13960 + 4)$ | 101 s | $e_1 = (\textbf{rerun } a \, b).$ | 1 |
| oneliners | spell | 1 | `iconv -f utf-8 -t ascii//translit` | 26404 $(= 12440 + 13960 + 4)$ | 59 s | $e_1 = (\textbf{concat } a \, b),$ <br> $e_2 = (\textbf{rerun } a \, b).$ | 2 |
| poets | sort_-words_-by_-rhyming | 6 | `rev` | 26404 $(= 12440 + 13960 + 4)$ | 59 s | $e_1 = (\textbf{concat } a \, b).$ | 1 |
| poets | count_-words | 1 | `IN=${IN:-../benchmarks/pipelines/poets/input/pg/}` <br> `sed "s;^;$IN;"` | 26404 $(= 12440 + 13960 + 4)$ | 59 s | $e_1 = (\textbf{concat } a \, b).$ | 1 |
| analytics-mts | vehicles per day | 1 | `sed 's/T..:..:..//'` | 26404 $(= 12440 + 13960 + 4)$ | 59 s | $e_1 = (\textbf{concat } a \, b).$ | 1 |
| analytics-mts | vehicle hours on road | 1 | `sed 's/T\(..\):..:../,\1/'` | 110444 $(= 59048 + 51392 + 4)$ | 126 s | $e_1 = (\textbf{concat } a \, b),$ <br> $e_2 = (\textbf{rerun } a \, b).$ | 2 |
| oneliners | top-n | 6 | `sed 100q` | 26404 $(= 12440 + 13960 + 4)$ | 258 s | $e_1 = (\textbf{rerun } a \, b).$ | 1 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| poets | trigram_-rec | 13 | `sed 5q` | 26404 (= 12440 + 13960 + 4) | 183 s | $e_1 = ($ **rerun** $a\ b)$. | 1 |
| unix50 | 7.3: decades unix re-leased | 5 | `sed s/\$/'0s'/` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = ($ **concat** $a\ b)$. | 1 |
| analytics-mts | vehicles per day | 5 | `sort` | 26404 (= 12440 + 13960 + 4) | 65 s | $e_1 = ($ **merge** $a\ b)$, $e_2 = ($ **merge** $b\ a)$, $e_3 = ($ **rerun** $a\ b)$, $e_4 = ($ **rerun** $b\ a)$. | 4 |
| poets | sort_-words_-by_-folding | 6 | `sort -f` | 26404 (= 12440 + 13960 + 4) | 65 s | $e_1 = ($ **merge**('-f') $a\ b)$, $e_2 = ($ **merge**('-f') $b\ a)$, $e_3 = ($ **rerun** $a\ b)$, $e_4 = ($ **rerun** $b\ a)$. | 4 |
| analytics-mts | vehicle days on road | 7 | `sort -k1n` | 26404 (= 12440 + 13960 + 4) | 65 s | $e_1 = ($ **merge**('-k1n') $a\ b)$, $e_2 = ($ **merge**('-k1n') $b\ a)$, $e_3 = ($ **rerun** $a\ b)$, $e_4 = ($ **rerun** $b\ a)$. | 4 |
| oneliners | shortest-scripts | 6 | `sort -n` | 26404 (= 12440 + 13960 + 4) | 64 s | $e_1 = ($ **merge**('-n') $a\ b)$, $e_2 = ($ **merge**('-n') $b\ a)$, $e_3 = ($ **rerun** $a\ b)$, $e_4 = ($ **rerun** $b\ a)$. | 4 |
| poets | sort | 6 | `sort -nr` | 26404 (= 12440 + 13960 + 4) | 65 s | $e_1 = ($ **merge**('-nr') $a\ b)$, $e_2 = ($ **merge**('-nr') $b\ a)$, $e_3 = ($ **rerun** $a\ b)$, $e_4 = ($ **rerun** $b\ a)$. | 4 |
| oneliners | sort-sort | 3 | `sort -r` | 26404 (= 12440 + 13960 + 4) | 66 s | $e_1 = ($ **merge**('-r') $a\ b)$, $e_2 = ($ **merge**('-r') $b\ a)$, $e_3 = ($ **rerun** $a\ b)$, $e_4 = ($ **rerun** $b\ a)$. | 4 |
| oneliners | top-n | 5 | `sort -rn` | 26404 (= 12440 + 13960 + 4) | 65 s | $e_1 = ($ **merge**('-rn') $a\ b)$, $e_2 = ($ **merge**('-rn') $b\ a)$, $e_3 = ($ **rerun** $a\ b)$, $e_4 = ($ **rerun** $b\ a)$. | 4 |
| analytics-mts | vehicles per day | 3 | `sort -u` | 26404 (= 12440 + 13960 + 4) | 64 s | $e_1 = ($ **merge**('-u') $a\ b)$, $e_2 = ($ **merge**('-u') $b\ a)$, $e_3 = ($ **rerun** $a\ b)$, $e_4 = ($ **rerun** $b\ a)$. | 4 |
| unix50 | 4.6: piece used most | 9 | `tail -n 1` | 26404 (= 12440 + 13960 + 4) | 104 s | $e_1 = ($ **first** $b\ a)$, $e_2 = ($ **second** $a\ b)$, $e_3 = (($ **back** '\n' **first**) $b\ a)$, $e_4 = (($ **fuse** '\n' **first**) $b\ a)$, $e_5 = (($ **back** '\n' **second**) $a\ b)$, $e_6 = (($ **fuse** '\n' **second**) $a\ b)$, $e_7 = ($ **rerun** $a\ b)$. | 7 |
| unix50 | 4.4: his-togram by piece | 1 | `tr ' ' '\n'` | 2700 (= 968 + 1728 + 4) | 41 s | $e_1 = ($ **concat** $a\ b)$, $e_2 = ($ **rerun** $a\ b)$. | 2 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| unix50 | 10.3: extract user-name | 7 | `tr '[A-Z]' '[a-z]'` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| unix50 | 4.5: his-togram by piece and pawn | 6 | `tr '[a-z]' 'P'` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| poets | merge_-upper | 3 | `tr '[a-z]' '[A-Z]'` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| unix50 | 9.1: extract word PORT | 3 | `tr '[a-z]' '\n'` | 26404 (= 12440 + 13960 + 4) | 59 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| poets | count_-vowel_-seq | 3 | `tr 'a-z' '[A-Z]'` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| unix50 | 8.4: longest words w/o hy-phens | 1 | `tr -c "[a-z][A-Z]" '\n'` | 2700 (= 968 + 1728 + 4) | 40 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| unix50 | 9.6: follow direc-tions | 9 | `tr -c '[A-Z]' '\n'` | 2700 (= 968 + 1728 + 4) | 40 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| unix50 | 9.8: TELE-communications | 1 | `tr -c '[a-z][A-Z]' '\n'` | 2700 (= 968 + 1728 + 4) | 40 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| oneliners | bi-grams | 1 | `tr -cs A-Za-z '\n'` | 2700 (= 968 + 1728 + 4) | 40 s | $e_1 = (\textbf{rerun } a\ b)$. | 1 |
| unix50 | 2.1: all Unix utilities | 2 | `tr -d ','` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| oneliners | spell | 5 | `tr -d '[:punct:]'` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| unix50 | 9.1: extract word PORT | 5 | `tr -d '\n'` | 2700 (= 968 + 1728 + 4) | 40 s | $e_1 = (\textbf{concat } a\ b)$, $e_2 = (\textbf{rerun } a\ b)$. | 2 |
| unix50 | 7.2: most fre-quent ma-chine | 6 | `tr -s ' ' '\n'` | 2700 (= 968 + 1728 + 4) | 41 s | $e_1 = (\textbf{rerun } a\ b)$. | 1 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| poets | count_-vowel_-seq | 4 | `tr -sc 'AEIOU' '[\012*]'` | 2700 (= 968 + 1728 + 4) | 40 s | $e_1 = (\textbf{rerun } a\, b)$. | 1 |
| poets | vowel_-sequencies_-gr_1K | 4 | `tr -sc 'AEIOUaeiou' '[\012*]'` | 2700 (= 968 + 1728 + 4) | 40 s | $e_1 = (\textbf{rerun } a\, b)$. | 1 |
| poets | count_-consonant_-seq | 4 | `tr -sc 'BCDFGHJKLMNPQRSTVWXYZ' '[\01 2*]'` | 2700 (= 968 + 1728 + 4) | 41 s | $e_1 = (\textbf{rerun } a\, b)$. | 1 |
| poets | merge_-upper | 4 | `tr -sc '[A-Z]' '[\012*]'` | 2700 (= 968 + 1728 + 4) | 40 s | $e_1 = (\textbf{rerun } a\, b)$. | 1 |
| poets | 2syllable_-words | 3 | `tr -sc '[A-Z][a-z]' ' [\012*]'` | 2700 (= 968 + 1728 + 4) | 40 s | $e_1 = (\textbf{rerun } a\, b)$. | 1 |
| poets | count_-words | 3 | `tr -sc '[A-Z][a-z]' '[\012*]'` | 2700 (= 968 + 1728 + 4) | 40 s | $e_1 = (\textbf{rerun } a\, b)$. | 1 |
| poets | sort_-words_-by_n_-syllables | 5 | `tr -sc '[AEIOUaeiou\012]' ' '` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat } a\, b)$, $e_2 = (\textbf{rerun } a\, b)$. | 2 |
| oneliners | bi-grams | 2 | `tr A-Z a-z` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat } a\, b)$, $e_2 = (\textbf{rerun } a\, b)$. | 2 |
| oneliners | diff | 2 | `tr [:lower:] [:upper:]` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat } a\, b)$, $e_2 = (\textbf{rerun } a\, b)$. | 2 |
| oneliners | diff | 5 | `tr [:upper:] [:lower:]` | 26404 (= 12440 + 13960 + 4) | 59 s | $e_1 = (\textbf{concat } a\, b)$, $e_2 = (\textbf{rerun } a\, b)$. | 2 |
| oneliners | bi-grams | 5 | `uniq` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = ((\textbf{stitch first})\, a\, b)$, $e_2 = ((\textbf{stitch second})\, a\, b)$, $e_3 = (\textbf{rerun } a\, b)$. | 3 |
| analytics-mts | vehicles per day | 6 | `uniq -c` | 26404 (= 12440 + 13960 + 4) | 59 s | $e_1 = ((\textbf{stitch2 `` add first})\, a\, b)$, $e_2 = ((\textbf{stitch2 `` add second})\, a\, b)$. | 2 |
| unix50 | 7.1: number of ver-sions | 3 | `wc -l` | 2700 (= 968 + 1728 + 4) | 39 s | $e_1 = ((\textbf{back `\textbackslash n` add})\, a\, b)$, $e_2 = ((\textbf{back `\textbackslash n` add})\, b\, a)$. | 2 |
| oneliners | shortest-scripts | 4 | `xargs -L 1 wc -l` | 26404 (= 12440 + 13960 + 4) | 103 s | $e_1 = (\textbf{concat } a\, b)$, $e_2 = ((\textbf{offset `` first})\, a\, b)$, $e_3 = ((\textbf{offset `` second})\, a\, b)$. | 3 |
| poets | count_-words | 2 | `xargs cat` | 26404 (= 12440 + 13960 + 4) | 60 s | $e_1 = (\textbf{concat } a\, b)$, $e_2 = ((\textbf{offset `` first})\, a\, b)$, $e_3 = ((\textbf{offset `` second})\, a\, b)$. | 3 |
| oneliners | shortest-scripts | 1 | `xargs file` | 26404 (= 12440 + 13960 + 4) | 61 s | $e_1 = (\textbf{concat } a\, b)$, $e_2 = ((\textbf{offset `` second})\, a\, b)$. | 2 |

**Table 7.** Performance results for benchmark scripts ($u_1 \geq 3\,\text{min}$)

| Benchmark | Script Name | Parallelized | Eliminated | $T_{\text{orig}}$ | $u_1$ | $u_{16}$ | $T_{16}$ |
|---|---|---|---|---|---|---|---|
| analytics-mts | 1.sh (vehicles per day) | 7/7 (7/7) | 3 (3) | 333 s (1.1×) | 376 s | 40 s (9.4×) | 29 s (13.1×) |
| analytics-mts | 2.sh (vehicle days on road) | 8/8 (8/8) | 3 (3) | 335 s (1.1×) | 379 s | 41 s (9.3×) | 28 s (13.5×) |
| analytics-mts | 3.sh (vehicle hours on road) | 8/8 (8/8) | 3 (3) | 408 s (1.0×) | 427 s | 51 s (8.4×) | 38 s (11.3×) |
| oneliners | bi-grams.sh | 3/5 (3/5) | 0 (0) | 668 s (1.5×) | 1007 s | 118 s (8.6×) | 115 s (8.7×) |
| oneliners | diff.sh | 4/7 (0/1, 2/2, 2/2, 0/1, 0/1) | 2 (0, 1, 1, 0, 0) | 325 s (1.5×) | 478 s | 98 s (4.9×) | 83 s (5.8×) |
| oneliners | nfa-regex.sh | 2/2 (2/2) | 1 (1) | 389 s (1.0×) | 391 s | 26 s (14.9×) | 27 s (14.7×) |
| oneliners | set-diff.sh | 5/8 (0/1, 3/3, 2/2, 0/1, 0/1) | 3 (0, 2, 1, 0, 0) | 879 s (1.5×) | 1308 s | 144 s (9.1×) | 128 s (10.2×) |
| oneliners | sort.sh | 1/1 (1/1) | 0 (0) | 273 s (1.4×) | 389 s | 39 s (10.0×) | 38 s (10.3×) |
| oneliners | spell.sh | 6/8 (6/8) | 3 (3) | 427 s (1.7×) | 736 s | 78 s (9.5×) | 61 s (12.1×) |
| oneliners | top-n.sh | 4/6 (4/6) | 1 (1) | 372 s (1.7×) | 622 s | 63 s (9.9×) | 50 s (12.4×) |
| oneliners | wf.sh | 4/5 (4/5) | 1 (1) | 1155 s (1.8×) | 2089 s | 196 s (10.7×) | 145 s (14.4×) |
| poets | 1_1.sh (count_words) | 4/6 (4/6) | 1 (1) | 360 s (1.8×) | 637 s | 84 s (7.6×) | 83 s (7.6×) |
| poets | 2_1.sh (merge_upper) | 5/7 (5/7) | 2 (2) | 307 s (1.8×) | 547 s | 79 s (6.9×) | 78 s (7.0×) |
| poets | 3_1.sh (sort) | 5/7 (5/7) | 1 (1) | 391 s (1.7×) | 665 s | 89 s (7.4×) | 88 s (7.6×) |
| poets | 3_2.sh (sort_words_by_folding) | 5/7 (5/7) | 1 (1) | 402 s (1.7×) | 681 s | 94 s (7.3×) | 94 s (7.2×) |
| poets | 3_3.sh (sort_words_by_rhyming) | 7/9 (7/9) | 2 (2) | 415 s (1.7×) | 699 s | 100 s (7.0×) | 100 s (7.0×) |
| poets | 4_3.sh (bigrams) | 4/8 (2/4, 0/1, 2/3) | 1 (1, 0, 0) | 635 s (1.4×) | 915 s | 173 s (5.3×) | 173 s (5.3×) |
| poets | 4_3b.sh (count_trigrams) | 4/9 (2/4, 0/1, 0/1, 2/3) | 1 (1, 0, 0, 0) | 862 s (1.2×) | 1049 s | 275 s (3.8×) | 279 s (3.8×) |
| poets | 6_1_2.sh (uppercase_by_type) | 4/6 (4/6) | 1 (1) | 330 s (1.9×) | 635 s | 64 s (10.0×) | 24 s (26.9×) |
| poets | 6_2.sh (4letter_words) | 7/11 (3/5, 4/6) | 2 (1, 1) | 327 s (2.0×) | 647 s | 80 s (8.1×) | 34 s (18.8×) |
| poets | 6_3.sh (words_no_vowels) | 5/7 (5/7) | 2 (2) | 220 s (1.1×) | 235 s | 32 s (7.4×) | 31 s (7.7×) |
| poets | 6_4.sh (1syllable_words) | 5/8 (5/8) | 2 (2) | 433 s (1.3×) | 542 s | 57 s (9.5×) | 31 s (17.4×) |
| poets | 6_5.sh (2syllable_words) | 5/8 (5/8) | 2 (2) | 397 s (1.1×) | 443 s | 48 s (9.2×) | 40 s (11.0×) |
| poets | 7_2.sh (count_consonant_seq) | 5/7 (5/7) | 2 (2) | 475 s (1.4×) | 678 s | 80 s (8.5×) | 48 s (14.2×) |
| poets | 8.2_1.sh (vowel_sequences_gr_1K) | 5/8 (5/8) | 1 (1) | 417 s (1.4×) | 573 s | 73 s (7.9×) | 42 s (13.7×) |
| poets | 8.2_2.sh (bigrams_appear_twice) | 4/9 (2/4, 0/1, 2/3, 0/1) | 1 (1, 0, 0, 0) | 645 s (1.4×) | 921 s | 177 s (5.2×) | 91 s (10.2×) |
| poets | 8.3_2.sh (find_anagrams) | 7/9 (2/4, 1/1, 1/1, 3/3) | 1 (1, 0, 0, 0) | 237 s (3.1×) | 724 s | 102 s (7.1×) | 50 s (14.5×) |
| poets | 8.3_3.sh (compare_exodus_genesis) | 6/10 (3/5, 1/2, 2/3) | 1 (1, 0, 0) | 334 s (2.0×) | 656 s | 74 s (8.8×) | 34 s (19.3×) |
| poets | 8_1.sh (sort_words_by_n_syllables) | 6/10 (3/5, 2/2, 1/3) | 2 (1, 1, 0) | 346 s (1.9×) | 653 s | 69 s (9.5×) | 26 s (24.6×) |
| unix50 | 14.sh (6.1: order bodies) | 3/3 (3/3) | 1 (1) | 143 s (1.3×) | 185 s | 31 s (6.0×) | 25 s (7.5×) |
| unix50 | 21.sh (8.4: longest words w/o hyphens) | 3/3 (3/3) | 1 (1) | 428 s (1.7×) | 733 s | 64 s (11.4×) | 49 s (14.9×) |
| unix50 | 23.sh (9.1: extract word PORT) | 6/6 (6/6) | 4 (4) | 111 s (1.8×) | 202 s | 23 s (8.8×) | 10 s (19.8×) |
| unix50 | 28.sh (9.6: follow directions) | 6/10 (6/10) | 3 (3) | 87 s (2.2×) | 188 s | 54 s (3.5×) | 49 s (3.8×) |
| **Total** | | 163/233 | 55 | | | | |
| **Max** | | | | 1155 s (3.1×) | 2089 s | 275 s (14.9×) | 279 s (26.9×) |
| **Min** | | | | 87 s (1.0×) | 185 s | 23 s (3.5×) | 10 s (3.8×) |
| **Mean** | | | | 420 s (1.6×) | 649 s | 85 s (8.2×) | 67 s (12.0×) |
| **Median** | | | | 389 s (1.5×) | 637 s | 74 s (8.5×) | 49 s (11.3×) |

**Table 8.** Combiners synthesized for all benchmark scripts

| Count | Synthesized Plausible Combiner |
|---|---|
| 81 | (**concat** *a b*) |
| 22 | (**rerun** *a b*) |
| 16 | (**merge**($*$) *a b*) or (**merge**($*$) *b a*) |
| 12 | ((**back** '\n' **add**) *a b*) or ((**back** '\n' **add**) *b a*) |
| 8 | (**rerun** *b a*) |
| 2 | ((**back** '\n' **first**) *a b*) or ((**back** '\n' **second**) *b a*) |
| 2 | (**first** *a b*) or (**second** *b a*) |
| 2 | ((**fuse** '\n' **first**) *a b*) or ((**fuse** '\n' **second**) *b a*) |
| 2 | ((**back** '\n' **second**) *a b*) or ((**back** '\n' **first**) *b a*) |
| 2 | (**second** *a b*) or (**first** *b a*) |
| 2 | ((**fuse** '\n' **second**) *a b*) or ((**fuse** '\n' **first**) *b a*) |
| 2 | ((**stitch2** ' ' **add first**) *a b*) or ((**stitch2** ' ' **add second**) *a b*) |
| 2 | ((**stitch first**) *a b*) or ((**stitch second**) *a b*) |

**Table 9.** Unsupported commands in all benchmark scripts

| Command | Reason Unsupported | Counterexample Input Streams |
|---|---|---|
| `awk "\$1 == 2 print \$2, \$3"` | KumQuat did not generate inputs for the command to produce nonempty outputs. | |
| `sed 1d` | **No combiners** $g$ **exist** such that $f(x_1 \mathbin{++} x_2) = g(f(x_1), f(x_2))$ for all streams $x_1, x_2$. | Each of $x_1, x_2$ has at least one line. |
| `sed 2d` | **No combiners** $g$ **exist** such that $f(x_1 \mathbin{++} x_2) = g(f(x_1), f(x_2))$ for all streams $x_1, x_2$. | Each of $x_1, x_2$ has at least two lines. |
| `sed 3d` | **No combiners** $g$ **exist** such that $f(x_1 \mathbin{++} x_2) = g(f(x_1), f(x_2))$ for all streams $x_1, x_2$. | Each of $x_1, x_2$ has at least three lines. |
| `sed 4d` | **No combiners** $g$ **exist** such that $f(x_1 \mathbin{++} x_2) = g(f(x_1), f(x_2))$ for all streams $x_1, x_2$. | Each of $x_1, x_2$ has at least four lines. |
| `sed 5d` | **No combiners** $g$ **exist** such that $f(x_1 \mathbin{++} x_2) = g(f(x_1), f(x_2))$ for all streams $x_1, x_2$. | Each of $x_1, x_2$ has at least five lines. |
| `tail +2` | **No combiners** $g$ **exist** such that $f(x_1 \mathbin{++} x_2) = g(f(x_1), f(x_2))$ for all streams $x_1, x_2$. | Each of $x_1, x_2$ has at least one line. |
| `tail +3` | **No combiners** $g$ **exist** such that $f(x_1 \mathbin{++} x_2) = g(f(x_1), f(x_2))$ for all streams $x_1, x_2$. | Each of $x_1, x_2$ has at least two lines. |